

ABSTRACT

Title of dissertation: DATA DRIVEN PREDICTION
WITHOUT A MODEL

Erin Michelle Lynch
Doctor of Philosophy, 2019

Dissertation directed by: Professor Eugenia Kalnay
Department of Atmospheric and
Oceanic Science
Dr. A. Surjalal Sharma
Department of Astronomy

Ensemble data assimilation techniques, including the Ensemble Transform Kalman Filter (ETKF), have been successfully used to improve prediction skill in numerical models for weather forecasting. However, less research has been conducted on data assimilation techniques for systems with no numerical model. In this study, we begin by applying the technique of bred vectors to a reconstructed phase space model for simple, autonomous nonlinear systems and compare the predictive capabilities of data driven bred vectors to those computed using a numerical model. Next, we show that a combination of the phase space reconstruction with ETKF yields a new technique, which we call Nearest Neighbor ETKF, for forecasting using only time series data. This technique is applied to a simple nonlinear system, the Lorenz three variable model, to demonstrate its effectiveness in forecasting noisy time series data. Finally, we use this technique to forecast field variations in the magnetosphere, which exhibit low dimensional behavior on the substorm time scale.

The time series data of the magnetic field variations monitored by the network of ground-based magnetometers in the auroral region are used for forecasting at two stages. In the first stage, the auroral electrojet indices, computed from the magnetometer data, are used to reconstruct the dynamics and Nearest Neighbor ETKF yields forecasts of the index that are more skillful than persistence. In the second stage, the multivariate time series from several auroral region magnetometers is used to reconstruct the phase space of the magnetosphere-solar wind system using Multi-channel Singular Spectrum Analysis. The Nearest Neighbor ETKF is applied to ensemble forecasts made using model data, constructed from long time series of the data from each magnetometer, in addition to observations in the reconstructed phase space, constructed from magnetometer measurements concurrent with the start of the forecast. Additionally, the spreads of the ensembles constructed to forecast these times series are used as precursors to understand and predict extreme space weather events.

DATA DRIVEN PREDICTION WITHOUT A MODEL

by

Erin Michelle Lynch

Dissertation submitted to the Faculty of the Graduate School of the
University of Maryland, College Park in partial fulfillment
of the requirements for the degree of
Doctor of Philosophy
April, 2019

Advisory Committee:

Professor Eugenia Kalnay, Co-Chair/Advisor

Dr. A. Surjalal Sharma, Co-Chair/Advisor

Professor James Carton

Professor Rachel Pinker

Professor Rajarshi Roy (Dean's Representative)

© Copyright by
Erin Michelle Lynch
2019

Acknowledgments

I owe a deep debt of gratitude to my advisors Dr. Eugenia Kalnay and Dr. Surja Sharma, without whom this definitely would not have been possible. I thank them both for their unwavering encouragement and support, their guidance, and for giving me faith in myself and my work.

I would like to thank my committee members for serving on both my dissertation committee and my prospectus committee and for their invaluable suggestions and feedback. I am also indebted to Prof. Kayo Ide for bringing me to the PhD program with AOSC, her instrumental contributions to this project, and guiding through my Master's degree. I thank the members of the Weather and Chaos group for providing me with the chance to present my work along the way and the insight and knowledge I've gained from their work.

I am grateful to Dr. Tony Mannucci, the members of the MRTIF group, and the NASA Living with a Star program as a whole, for the financial support for this work as well as the truly invaluable experience and feedback I was able to receive as part of this project.

I would like to thank Dr. Xi Shao who helped me to get my current job working at NOAA, really a dream come true. I am extremely grateful to Dr. Changyong Cao for his mentorship and guidance in building the next stage of my career.

Most especially, I would like to thank my family and friends: my mother for more than I could possibly list here; my siblings for being constant sources of inspiration; my grandmother and late grandfather for all their endless support and

confidence in me; my friends Cathy, Rebekah, and Clare for showing me the way; Amanda for her friendship and help through the toughest of times; and Jo and Angie for being there from the start.

Table of Contents

Acknowledgements	ii
List of Tables	vi
List of Figures	viii
1 Introduction	1
1.1 Background and Motivation	1
1.2 Tools for Numerical Weather Prediction	2
1.3 Nonlinear Time Series Prediction	5
1.4 Forecasting Space Weather	6
1.5 Thesis Objectives	8
2 Data Driven Ensemble Transform Kalman Filter	10
2.1 Dynamical Modeling Using Time Series Data	10
2.2 Time Delay Embedding	11
2.3 Reconstructing Phase Space by Singular Spectrum Analysis	14
2.4 The Ensemble Transform Kalman Filter	16
2.5 Nearest Neighbor Ensemble Transform Kalman Filter	18
2.5.1 Model Construction	19
2.5.2 Forecast Cycle	19
2.5.3 Reconstruction of Original Time Series	21
2.5.4 Localization	22
2.5.5 Data Density	22
2.6 Multi-channel Extension	23
3 Bred Vectors in the Reconstructed Phase Space	24
3.1 Introduction	24
3.1.1 Method of Bred Vectors	26
3.1.2 Nearest-Neighbor Bred Vectors	27
3.1.3 Experimental Setup	28
3.2 Bred Vectors in the Reconstructed Phase Space of the Lorenz System	29
3.2.1 Reconstruction of the Lorenz Attractor	31

3.2.2	Bred Vectors Results	34
3.2.3	Predicting Regime Changes	36
3.3	Bred Vectors in the Reconstructed Phase Space of the Chua Oscillator	40
3.3.1	Reconstruction of the Chua Attractor	43
3.3.2	Bred Vectors Results	44
3.3.3	Predicting Regime Change	47
3.4	Bred Vectors in the Reconstructed Phase Space of the Rössler System	48
3.4.1	Reconstruction of the Rössler Attractor	50
3.4.2	Bred Vectors Results	51
3.4.3	Predicting Regime Change	53
3.5	Summary and Conclusions	56
4	Modeling and Forecasting the Lorenz System	59
4.1	Introduction	59
4.2	The Perfect Model Case	60
4.2.1	Results for the Lorenz System	62
4.2.2	Tuning Embedding Parameters	64
4.3	Additive Observational Noise	65
4.3.1	Phase Space Reconstruction of Noisy Data	65
4.4	Discussion and Conclusions	70
5	Modelling and Forecasting Space Weather Using Time Series Data of Magnetic Field Variations	72
5.1	Introduction	72
5.1.1	Geomagnetic Substorm Dynamics	72
5.1.2	Reconstruction of Magnetospheric Dynamics	76
5.2	Forecasting the AL Index During Geomagnetic Substorms	80
5.2.1	Experimental Setup	80
5.2.2	Forecasting the AL Index during High Speed Stream Events	82
5.2.3	Forecasting the AL Index during Coronal Mass Ejections	86
5.3	Forecasting Ground Based Magnetometer Measurements	89
5.4	Forecasting Extreme Events using Ensemble Spread	92
5.5	Conclusions	96
6	Summary and Conclusions	98
	Bibliography	103

List of Tables

3.1	Contingency tables based on the rule that regime change will occur in the orbit following the appearance of high growth rate bred vectors using three different methods. In (b) and (c) using the nearest-neighbor breeding, high growth rate points in orbits with absolute values of extrema above 1 are excluded. OBS and FCST stand for observed and forecast, respectively; (a)-(c) are the same as in Fig. 3.5.	38
3.2	Measures of forecast accuracy in terms of the Hit Rate (HR), Threat Score (TS), and False Alarm Rate (FAR); (a)-(c) are the same as in Fig. 3.5. The final row shows the values when the threshold of $x(t_i)$ rule is used.	39
3.3	Contingency tables based on the rule that regime change will occur in the orbit following the appearance of high growth rate bred vectors using three different methods. In (b) and (c) using the nearest-neighbor breeding, high growth rate points in orbits with absolute values of extrema above 1 are excluded. OBS and FCST stand for observed and forecast, respectively; (a)-(c) are the same as in Fig. 3.12.	48
3.4	Measures of forecast accuracy in terms of the Hit Rate (HR), Threat Score (TS), and False Alarm Rate (FAR); (a)-(c) are the same as in Fig. 3.12.	49
3.5	Contingency tables based on the rule that regime change will occur in the orbit following the appearance of high growth rate bred vectors using three different methods. In (b) and (c) using the nearest-neighbor breeding, high growth rate points in orbits with absolute values of extrema above 1 are excluded. OBS and FCST stand for observed and forecast, respectively; (a)-(c) are the same as in Fig. 3.18.	56
3.6	Measures of forecast accuracy in terms of the Hit Rate (HR), Threat Score (TS), and False Alarm Rate (FAR); (a)-(c) are the same as in Fig. 3.18. The final row shows the values when the threshold of $x(t_i)$ rule is used.	57

5.1	Normalized root mean squared errors for forecasts of various HSS events. Forecasts are made using NN ETKF and persistence, predicting the value of the AL index 20, 40, and 60 minutes beyond the current observation.	85
5.2	Normalized root mean squared errors for forecasts of various CME events. Forecasts are made using NN ETKF and persistence, predicting the value of the AL index 20, 40, and 60 minutes beyond the current observation.	88

List of Figures

3.1	The attractor of the Lorenz 3 variable system in the natural 3 dimensional phase space of (x, y, z)	30
3.2	The time series data for the variables $x(t)$, $y(t)$, and $z(t)$ for the Lorenz system. The signature of regime change can be seen in the $x(t)$ variable when a sign change occurs.	31
3.3	3.3(a) The mutual information function as a function of time delay for the Lorenz $x(t)$ variable. 3.3(b) The percentage of false nearest neighbors for attractors reconstructed from the $x(t)$ variable of the Lorenz system for various time delays as a function of embedding dimension. For a range of time delays, the percentage of false nearest neighbors becomes negligible for more than three dimensions.	32
3.4	The reconstructed attractor in the embedded phase space formed by time delay vectors with dimension $m = 3$ and time delay $\tau = 7$	33
3.5	Growth rates of bred vectors in the Lorenz system using three different methods: 3.5(a) standard breeding using the ordinary differential equations in the phase space (x, y, z) ; 3.5(b) nearest-neighbor breeding in the phase space (x, y, z) ; and 3.5(c) nearest-neighbor breeding in the reconstructed phase space (x_1, x_2, x_3) . The colored points correspond to negative (blue), low (green), medium (yellow) and high (red) growth; see text for the value of the thresholds.	34
3.6	The first coordinate of phase space as a function of time, with red stars indicating the points with high growth rate ($g_i \geq 6.4$) bred vectors; (a)-(c) are the same as in Fig. 3.5.	35
3.7	The observed number of cycles vs the number of high growth rate bred vectors observed in previous regime. The numbers indicate the frequency with which a given pair was observed. The line correspond to the values in prediction rule 2. (a) Bred vectors are computed using the ODEs. (b) NN bred vectors are computed using the full phase space model. (c) NN bred vectors are computed using the reconstructed phase space model.	40
3.8	The attractor of the Chua 3 variable system in the natural 3 dimensional phase space of (x, y, z)	41

3.9	The time series data for the variables $x(t)$, $y(t)$, and $z(t)$ for the Chua system. The signature of regime change can be seen in the $x(t)$ variable when a sign change occurs.	42
3.10	3.10(a) The mutual information function as a function of time delay for the Chua $x(t)$ variable. 3.10(b) The percentage of false nearest neighbors for attractors reconstructed from the $x(t)$ variable of the Chua system for various time delays as a function of embedding dimension. For a range of time delays, the percentage of false nearest neighbors becomes negligible for more than three dimensions.	43
3.11	The reconstructed Chua attractor in the embedded phase space of time delay vectors with $m = 3$ and $\tau = 16$	44
3.12	Growth rates of bred vectors in the Chua system using three different methods: 3.12(a) standard breeding using the ordinary differential equations in the phase space (x, y, z) ; 3.12(b) nearest-neighbor breeding in the phase space (x, y, z) ; and 3.12(c) nearest-neighbor breeding in the reconstructed phase space (x_1, x_2, x_3) . The colored points correspond to negative (blue), low (green), medium (yellow) and high (red) growth; see text for the value of the thresholds.	45
3.13	The first coordinate of phase space as a function of time, with red stars indicating the points with high growth rate ($g_i \geq 6.4$) bred vectors; (a)-(c) are the same as in Fig. 3.12(a).	46
3.14	The attractor of the Rössler 3 variable system in the natural 3 dimensional phase space of (x, y, z)	50
3.15	The time series data for the variables $x(t)$, $y(t)$, and $z(t)$ for the Rössler system. The signature of regime change can be seen in the $x(t)$ variable when a sign change occurs.	51
3.16	3.16(a) The mutual information function as a function of time delay for the Rössler $x(t)$ variable. 3.16(b) The percentage of false nearest neighbors for attractors reconstructed from the $x(t)$ variable of the Rössler system for various time delays as a function of embedding dimension. For a range of time delays, the percentage of false nearest neighbors becomes negligible for more than three dimensions.	52
3.17	The reconstructed attractor in the embedded phase space formed by time delay vectors with dimension $m = 3$ and time delay $\tau = 34$	52
3.18	Growth rates of bred vectors in the Lorenz system using three different methods: 3.18(a) standard breeding using the ordinary differential equations in the phase space (x, y, z) ; 3.18(b) nearest-neighbor breeding in the phase space (x, y, z) ; and 3.18(c) nearest-neighbor breeding in the reconstructed phase space (x_1, x_2, x_3) . The colored points correspond to negative (blue), low (green), medium (yellow) and high (red) growth; see text for the value of the thresholds.	54
3.19	The first coordinate of phase space as a function of time, with red stars indicating the points with high growth rate ($g_i \geq 6.4$) bred vectors; (a)-(c) are the same as in Fig. 3.18(a).	55

4.1	Comparison of forecasts made using three versions of the perfect model: 4.1(a) the ODE model, 4.1(b) the full phase space model (x, y, z) , and 4.1(c) the phase space model (x_1, x_2, x_3) reconstructed from $x(t)$	62
4.2	The average root mean squared error in both the analysis (blue) and the background (red) as a (a) function of embedding dimension, (b) time delay, and (c) the length of the forecast window.	63
4.3	(a)The mutual information function as a function of time delay and (b) the fraction of false nearest neighbors as a function of embedding dimension for the Lorenz $x(t)$ plus various levels of additive noise variable.	66
4.4	An enlarged portion of the eigenvalue spectrum computed for various levels of additive noise for the Lorenz $x(t)$ variable to show the noise floor.	67
4.5	The first six eigenvectors of the Lorenz system x variable with additive noise levels of 0% and 50%.	68
4.6	Comparison of reconstructed forecasts made of the Lorenz $x(t)$ variable with a noise level 60% of the standard deviation of the x variable. Panel (a) depicts the background and analysis of NN ETKF forecasts made in the full (x, y, z) phase space. Panel (b) depicts the background and analysis of NN ETKF performed in the reconstructed (x_1, x_2, x_3) phase space. Panel (c) depicts the background and analysis of NN ETKF forecasts in the space of the first six principal components after SSA was applied to the reconstructed phase space.	69
4.7	The root mean squared error in the mean values of the forecast and analysis ensembles, as a percentage of the standard deviation in the time series $x(t)$ vs. the percentage of error added to the time series to simulate observations.	70
5.1	Schematic of the interaction between the solar wind magnetic field and the Earth's magnetic field	74
5.2	The locations of the 12 ground-based magnetometer stations whose measurements contribute to the construction of the AE indices	75
5.3	The AL index (orange curve) is the lower envelope of anomalies measured at each of the stations (gray curves). Data from eight of the twelve stations that contribute to the AL index are shown during a high speed stream in May of 2011.	75
5.4	The mutual information function as a function of time delay and the fraction of false nearest neighbors as a function of embedding dimension for the AL index.	78
5.5	The percent variance explained by each mode of the eigenspectrum on the AL index. After 6 modes, nearly 90% of the variance in the signal is explained.	79
5.6	The first six eigenvectors and principal components for a substorm even in 2005.	79

5.7	Schematic of the co-rotating interaction region (CIR) that forms during an HSS event [1]	83
5.8	Forecasts made of the AL index during an HSS event that occurred in April of 2011. The green and red curves represent forecasts made using persistence and NN ETKF respectively. 5.8(a) depicts a series of 20 minute forecasts, 5.8(b) depicts a series of 40 minute forecasts, and 5.8(c) depicts a series of 60 minute forecasts. The skill score of the NN ETKF forecasts with respect to persistence is quoted in the title.	84
5.9	The correlation between the true value of the AL Index during the April 2011 HSS event and forecasts made using the NN ETKF and persistence.	86
5.10	Schematic the structure of a CME [2]	87
5.11	Forecasts made of the AL index during a CME event that occurred in September of 2011. The green and red curves represent forecasts made using persistence and NN ETKF respectively. 5.11(a) depicts a series of 20 minute forecasts, 5.11(b) depicts a series of 40 minute forecasts, and 5.11(c) depicts a series of 60 minute forecasts. The skill score of the NN ETKF forecasts with respect to persistence is quoted in the title.	89
5.12	Forty minute forecasts of the data from each of the 8 stations that had data available to the public and the skill scores of the forecasts with respect to persistence. Skill scores range from -0.05 to 0.27.	90
5.13	The left panel depicts the scenario in which an ensemble is not subject to an extreme event, while the right panel shows the effect of an extreme event on the spread of the ensemble members. Figure from [3].	93
5.14	The spread of the nearest neighbors to a control trajectory over a series of forecasts is indicated by colored stars along the $x(t)$ variable. This figure was provided by Keenan Eure.	94
5.15	The dashed line depicts the observed value of the AL index while the colored circles indicate the mean value of the ETKF forecast ensemble. The color and size of the circles correspond to the magnitude of the ensemble spread. The ensemble spread is well correlated with the value of the AL index.	95
5.16	The black curve represents the observed value of the AL index during the April 2011 substorm while the blue curve is a forecast made by the linear regression in Eq. 5.3.	95

Chapter 1: Introduction

1.1 Background and Motivation

Nonlinear systems are ubiquitous in nature. A key feature of nonlinear dynamics is an inherent difficulty in predicting the evolution of the system. Even with full knowledge of the underlying dynamics, small errors in initial conditions can grow exponentially over time causing large forecast errors. This phenomena is well known in the case of weather forecasting where capturing highly nonlinear behavior becomes challenging. Tools for mitigating these effects in numerical forecasts and improving predictability have been developed extensively.

However, many systems in nature are difficult or impossible to forecast using numerical models. Constructing a first principles model that can be solved numerically is challenging when physical processes are not well understood. In the case that a model can be constructed, discretization of such models require that physical processes that occur on smaller scales than those represented by discrete grid be parameterized. This can hamper the ability of a numerical model to capture the all the relevant behavior of the system. Finally, to initialize a model forecast, sufficiently many state variables must be observed to capture the current state. Often only a partial state vector or perhaps a scalar value can be measured. The predictability

of many systems is impaired by these limitations and more; however, forecasting such nonlinear systems from time series data has been an area of active research for some time [4–10].

The motivation for this work comes from a desire to combine techniques common to numerical weather prediction (NWP) that improve forecast skill with data driven techniques for predicting the evolution of nonlinear systems based on time series data. The potential applications for these new techniques are vast, but here they will be used to forecast a natural phenomenon that is particularly difficult to forecast numerically - the development of a class of space weather events called magnetospheric substorms driven by the solar wind.

1.2 Tools for Numerical Weather Prediction

The past century has seen remarkable advances in our ability to forecast complex, nonlinear systems. Today, NWP accomplishes a feat that seemed insurmountable at the turn of the previous century when the equations governing the dynamics of the atmosphere were solved by hand. [11–13]. Great effort over many decades has brought the state of forecasting the behavior of the neutral atmosphere to its current level of success.

While the first principles governing the flow of air within the atmosphere have been well known, their complicated and nonlinear nature has made direct solutions of such equations near impossible until the advent of the computer age. Computers can perform the many computations required to numerically solve the governing equa-

tions on increasingly fine grids. Since then, a number of advances have contributed to the current success seen in forecasting the weather. The first is that computers have become better over time. This allows for the inclusion of more physics and more refined solution domains. The second major advance is seen in the number of observations available. The earliest days of weather forecasting depended upon individuals making regular detailed measurements of their local conditions. Now, in addition to many in situ measurements taken by all manner of weather stations, airplanes, and sea-crafts in addition to more specialized vehicles like radiosondes, weather balloons, and unmanned submersibles, we enjoy global coverage from remote sensing satellites that provide radiances from which atmospheric variables can be retrieved through radiative transfer models. The third major area in which advances have enabled enhanced numerical weather prediction is in the development of tools to help capture these inherent instabilities and mitigate their effects.

One such tool is bred vectors [14,15]. Breeding takes advantage of the divergent nature of nearby trajectories in nonlinear systems to characterize these instabilities in a computationally efficient way. These vectors capture the fastest growing modes of variability and, since their inception, these vectors have been applied to testbeds of simple nonlinear system to further explore their properties. In fact, bred vectors themselves have been shown to provide predictive skill for certain regime transitions [16].

Even with a perfect model and increased observational coverage, it is still a challenge to initialize complex, nonlinear models such as those for NWP. Each of the model variables must be known at each location which is impossible even

with enhanced observational coverage. Additionally, even the best observations include errors. The nonlinear nature of weather models that made them so difficult to perfect in the first place ensure that errors in the initial conditions will grow exponentially during the forecast window, to potentially complete divergence from the true evolution of the atmosphere. Data assimilation is a tool that serves to construct initial conditions with improved errors relative to the errors introduced by observations and previous model forecasts. By optimally averaging between the two estimates available at the start of the forecast window, taking into account their relative errors, an estimate of the true state of the system can be obtained.

Within the field of data assimilation, advances over the years have allowed for improvements beyond the use of static background error covariances in 3D variational methods, such as the Kalman filter [17], to those that incorporate some aspect of time. One framework assimilate asynchronous observations using a variational method called 4D VAR [18, 19]. This requires the use of a numerical model because the adjoint must be obtained to handle observations at different times. On the other hand, ensemble methods evolve background error covariances in time by computing them based on multiple model realizations from a number of perturbed initial conditions, thus sampling the variability in the model space [20, 21]. This technique has great potential to be applied to data driven methods because the ensemble members need not necessarily come from numerical model forecasts. The particular variant of an ensemble Kalman filter used in this study is the Ensemble Transform Kalman Filter (ETKF) [22]. Recently, efforts to extend the capability and forecast improvements due to data assimilation to systems lacking a numerical

model have been made [23], including this study.

1.3 Nonlinear Time Series Prediction

Complicated models and a dearth of data necessitates the exploration of alternative forecasting strategies. Many nonlinear systems that are difficult to observe and model benefit from a data driven approach to forecasting. The study of various systems through measurements of a partial set of variables, even in some cases a single scalar variable, has yielded predictive skill in forecasting future behavior.

Given a set of historical data from a system there are many approaches that can be taken. Making parametric fits of curves with various functional forms to datasets works well if the solution follows a known form. However, when there is no closed form solution to the data, or a linear fit to a polynomial is not appropriate, more advanced techniques need to be explored. An extension of a linear model is an auto-regressive moving average models (ARMA) that includes a stochastic component to account for some degree of nonlinearity. On the other hand, neural networks and other deep learning techniques, where a hidden layers of interconnected “neurons” transform the in-put variables to the out-put variables, can recreate complex, nonlinear behavior without knowing anything about the form of the solution.

For time series data arising from a dynamical system, i.e. one whose evolution is defined by a phase space [24], we can reconstruct the behavior of the system, even if only a single variable is observed. Through time delay embedding, developed by Packard et al. [25] and formalized into a theorem by Takens [26], phase space

models that account for all of the dynamics of a system can be constructed given a sufficiently long time series of observations. Behavior of arbitrary initial conditions can be forecasted by finding analogs among the reconstructed phase space states [27].

Fortunately, many systems in nature that we are interested in forecasting are deterministic, perhaps with underlying dynamics that may not be known to us. If they are also dissipative, i.e. the volume of the phase space of the initial state contracts under the dynamics, then the phase space model can be low dimensional and this data driven forecasting technique becomes viable. Phase space models constructed in this way will become the substitute for those made by numerical models in the techniques developed for NWP.

1.4 Forecasting Space Weather

Space weather is a class of phenomena that have recently garnered much attention. While space weather effects have been visible in the polar region aurora, the potentially devastating effects of which are only just being realized. While the effects are largely present in the naturally occurring subsystems of the earth's atmosphere and magnetosphere, man-made objects are severely affected. As we build structures along the earth's surface susceptible to the devastating effects of geomagnetically induced currents and structures within the near earth space environment vulnerable to the effects of highly energetic ionized particles, we stand to lose vital technologies upon which we have become dependant. Potentially crippling fallout from such large scale events include the severe economic impacts of having to rebuild vital

infrastructure including power grids and railway lines, the loss of communication and other satellites upon which the ease of modern life depend, and the potential loss of life and property.

The need to forecast such events is evident. However the current state of space weather forecasting lags considerably behind that of its terrestrial counterpart. Modeling space weather suffers from many drawbacks. First, the scale of the system is vast covering many relevant characteristic time and length scales. To truly capture the system one must model from “the sun to the mud” covering both the evolution of solar events such as high speed streams, flares, and coronal mass ejections; their propagation through and interaction with the space environment as the solar wind; and finally their interaction with a number of coupled earth systems including the magnetosphere, ionosphere, neutral atmosphere and potentially man-made earth-based systems to truly capture the full extent of the effects. The efficient modeling of many of these systems is still an area of active research. Physics governing the evolution of solar structures and the interaction of magnetic fields within key areas of the magnetosphere are still not well understood.

The scale of the system also contributes to another drawback, a lack of observations. Currently, space weather forecasts begin when activity is detected on the sun. The sun is well observed along the line of sight between it and the earth by several space based observatories. However to determine the direction of material ejected from its surface, multiple views are needed. Even in this limited context, adequate observations are difficult to come by. Stereo A and B orbit in opposite directions, potentially providing two additional views from different angles allowing

for the determination of direction and extent of material ejected from the sun. In other regions of the system observations are even more scarce. Only two satellites make in situ observation of the solar wind plasma parameters and magnetic field currently. The ACE satellite sits at the first Lagrange point and measures the upstream solar wind before it reaches earth. More recently, an effort is underway to observe the magnetosphere at key locations. But largely, observations of the magnetosphere made at the surface of the earth by a number of ground based magnetometers are needed to complement the in situ spacecraft observations.

1.5 Thesis Objectives

The objective of this study is to understand the effect of data assimilation on forecasts made using data derived models. As a preliminary step, we also explore the predictive skill of the growth rate of bred vectors using data derived model. We constructed these models from time series observations of univariate or multivariate phase space variables. The approach we developed combines phase space reconstruction to create this model with the ensemble transform Kalman filter to make improved forecasts. To test this approach, we use the Lorenz three variable model. This allows us to compare directly the performance of forecasts made in the reconstructed phase space to forecasts made using the known dynamical equations.

Of particular interest to this study are time series that characterize space weather conditions. The signature of substorm development can be seen in the scalar auroral electrojet indices. These indices have been studied for decades and

are used throughout the literature to identify substorm events. We use our technique to forecast the value of the AL index to predict such events.

This technique can be extended to multivariate time series. Rather than reducing the observed anomalies in the auroral oval to a scalar value, we also forecast the magnetometer measurements simultaneously. Forecasts of the multivariate time series provide additional information including timing and location of anomalies. Finally, the ensemble forecasts made of the AL index are used to identify and predict the occurrence of extreme events.

This thesis is organized as follows: Chapter 2 describes the techniques combined to perform the study. A preliminary study of bred vectors in the reconstructed phase space is described in Chapter 3. Chapter 4 includes the results of the Nearest Neighbor Ensemble Transform Kalman Filter applied to the Lorenz 3 variable system. Chapter 5 contains the results of applying this technique to the lower auroral electrojet index and the results of the multi-variate extension to the auroral oval magnetometer data. Results utilizing the spread of the forecast ensembles to predict regime change in the Lorenz system as well as the magnitude of space weather events are also presented in Chapter 5. Finally, some conclusions are summarized in Chapter 6.

Chapter 2: Data Driven Ensemble Transform Kalman Filter

Forecasting nonlinear systems with no numerical model and limited observations requires that we take full advantage of the data that is available and exploit the nonlinearity of the system. First a data-derived model is constructed using nonlinear time series analysis techniques of time-delay embedding. Since observations are contaminated by high dimensional noise, singular spectrum analysis (SSA) is applied to the phase space model to filter out unwanted noise, capture the dynamical signal, and create an orthogonal basis for the phase space model. Finally, the Ensemble Transform Kalman Filter (ETKF) combines forecasts made using the phase space model with observations of the current state of the system to improve the skill of forecasts by improving the initial conditions of each subsequent forecast cycle.

2.1 Dynamical Modeling Using Time Series Data

Forecasting future values of a scalar, nonlinear time series of observations has been a rich field of study [4,7,8]. From making simple linear predictor models based on a training set of data, to improved models that incorporate nonlinear behavior, there are many techniques available to those interested in tackling such a problem

[24, 28, 29]. For nonlinear dissipative systems, the evolution follows trajectories that lie on an attractor. For an infinitely long time series, the system will continue to visit each neighborhood of the attractor infinitely many times, coming arbitrarily close to previous trajectories without ever intersecting. For dissipative systems, the dimension of the attractor itself is smaller than the dimension of the phase space in which it is contained because phase space volumes contract as the system evolves in time [30].

Given an attractor that is densely populated by a long phase space trajectory, predictions can be made using analogs of arbitrary initial conditions. Nearby trajectories will tend to follow a similar evolution in the short term [8]. In the long term, the nonlinear nature of the governing dynamics will cause nearby trajectories to diverge exponentially. Given a long historical time series, the attractor of the system can be covered such that analogs for arbitrary initial conditions can be found within neighborhoods of a given radius.

2.2 Time Delay Embedding

For dissipative, nonlinear, systems for which the full state vector cannot be observed, it is possible to reconstruct the leading features of the attractor from a time series of scalar observations for systems that contract onto a low dimensional attractor. The reconstructed attractor is diffeomorphic with the original attractor and preserves the original topology. [25, 26] Given a time series of data, $\{x(t_1), x(t_2), \dots, x(t_N)\}$, taken from measurements at regular intervals, multivariate

state vectors can be constructed by taking as each component a time-lagged value of the original series [24, 31].

$$\tilde{\mathbf{x}}(t_{i+(m-1)\tau}) = (x(t_i), x(t_{i+\tau}), \dots, x(t_{i+(m-1)\tau}))^T \quad (2.1)$$

The time delay between each successive component is denoted by τ , and should be long enough that within that interval, information about the unobserved, but non-linearly coupled, variables is incorporated in the evolution of the observed variable. If the time delay is too long the observations at the beginning and end of the interval will be uncoupled and therefore essentially independent. To ensure that each new component of the so called time-delay vector introduces the optimal amount of new information, the time delay can be determined by looking at the decay of the information shared by each pair of components. This quantity is measured by the average mutual information function [27]. The autocorrelation function can also be used to obtain the time delay, τ ; however, the mutual information function yields more reliable values because it can capture nonlinear relationships while the autocorrelation function depends on linear relationships.

The time delay is determined from the mutual information function, given by

$$I(\tau) = \sum_{t=0}^N p(x(t), x(t+\tau)) \log \left[\frac{p(x(t), x(t+\tau))}{p(x(t))p(x(t+\tau))} \right], \quad (2.2)$$

where $p(x(t))$ and $p(x(t+\tau))$ are the probability distributions of the variable $x(t)$ and the time lagged value respectively.

The mutual information function is a measure of the extent to which knowledge of one random variable provides knowledge of another random variable. If the two are mutually independent, they share no common information and knowledge of one does not provide knowledge of the other, thus the mutual information function tends to zero. It is similar to an auto-correlation in that it will pick up linear relationships between the two lagged variables, but it is also possible to detect nonlinear relationships.

The time delayed reconstruction depends on information about unobserved variables being introduced by each new component of the time delayed vector. The time delay selected should be long enough that any two components do not provide redundant information, but not so long that the nonlinear nature of the system has eliminated any correlation between the two

The dimension of the state vector, denoted by m , should be large enough that the reconstructed attractor is completely unfolded and trajectories lying on it are proper phase space trajectories, i.e. they do not intersect. If the dimension of the embedding is too small, phase space trajectories will necessarily intersect. Thus, points that fall on these intersecting trajectories will appear to be nearest neighbors when they should not be nearby at all. By estimating the percentage of "false nearest neighbors", those points that are no longer neighbors to one another when the dimension of the embedding increases, one can determine the optimal embedding dimension by watching this quantity fall to zero [6, 32, 33].

2.3 Reconstructing Phase Space by Singular Spectrum Analysis

There are several considerations that must be made for dealing with real data. The first is the duration of the time series. It is important that there be sufficiently many observations to cover the attractor densely enough to find suitable nearest neighbor analogs [24].

The second concern is that real data often contains random noise in addition to the underlying signal of interest. This noise is high dimensional and can increase the estimated dimension of the attractor. We wish to filter out this noise by using Singular Spectrum Analysis (SSA) [9, 34]. Given the m -dimensional reconstructed attractor, SSA allows for the identification of the directions that contain most of the variability seen in the data. The remaining dimensions that contain minimal variability correspond to the noise present in the data.

The application of SSA also serves to construct a proper basis for the model. The basis of the time delay vectors is not orthogonal. The basis of the principal components constructed using SSA are eigenvectors and are therefore orthogonal to one another. As analogs for our forecast initial conditions, we use nearest neighboring points among the model trajectory located by minimizing the distance to the arbitrary initial condition. Thus, with a proper basis, the distances to nearest neighbors in the space of principal components is well defined.

To perform the SSA, we first construct the trajectory matrix of time delay

vectors. From the trajectory matrix,

$$\mathbf{D} = \begin{pmatrix} x(t_1) & x(t_2) & \cdots & x(t_{N'}) \\ x(t_{1+\tau}) & x(t_{2+\tau}) & \cdots & x(t_{N'+\tau}) \\ \vdots & \vdots & \ddots & \vdots \\ x(t_{1+(m-1)\tau}) & x(t_{2+(m-1)\tau}) & \cdots & x(t_{N'+(m-1)\tau}) \end{pmatrix} \quad (2.3)$$

the covariance matrix is given by

$$\mathbf{C} = \frac{1}{N} \mathbf{D} \mathbf{D}^T. \quad (2.4)$$

From the covariance matrix, eigenvalue decomposition is used to find the series of eigenvalues, $\{\lambda_1, \lambda_2, \dots, \lambda_m\}$, and eigenvectors $\{\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_m\} \in \mathbb{R}^m$. The eigenvalues correspond to the amount of variance attributable to the mode defined by the corresponding eigenvector. The eigenvalues tend to fall off with increasing mode number. The eigenvectors that correspond to eigenvalues that fall below a certain noise floor can be discarded. Taking only the first k modes corresponding to the largest eigenvalues, an $m \times k$ transformation matrix given by

$$\mathbf{T} = [\mathbf{e}_1 \mathbf{e}_2 \cdots \mathbf{e}_k]^T \quad (2.5)$$

can be used to map the m dimensional time delay state vectors onto the k dimensional attractor that describes the dynamical signal of the data, yielding a k dimensional phase space of the principal component time series.

2.4 The Ensemble Transform Kalman Filter

Both models used to predict future states of a system and observations used to initialize the models contain errors. For nonlinear systems these errors compound over time, leading to a rapid divergence between the forecast and the true state of the system being predicted. Data assimilation is a technique used to systematically improve the initial conditions used to initiate a forecast by taking into account both the errors in the observations and the errors introduced by the previous model forecast to make a better estimate of the true state of the system to start the next forecast. [13]

The improved estimation of the true state is called the analysis given by \mathbf{x}^a , while the a priori estimate of the state provided by the previous model forecast is called the background and denoted by \mathbf{x}^b . To find the analysis given the background state and an observation \mathbf{x}^o , minimize the cost function that assumes both the background and observation have Gaussian error statistics with respect to the true state of the system:

$$J(\mathbf{x}) = \frac{1}{2}((\mathbf{x} - \bar{\mathbf{x}}^b)^T(\mathbf{P}^b)^{-1})(\mathbf{x} - \bar{\mathbf{x}}^b) + [\mathbf{y}^o - \mathbf{H}^o\mathbf{x}]^T \mathbf{R}^{-1} [\mathbf{y}^o - \mathbf{H}^o\mathbf{x}] \quad (2.6)$$

where \mathbf{P}^b and \mathbf{R} are the error covariances of the background and observations respectively. The observation operator \mathbf{H}^o takes vector \mathbf{x} from the state space of the system and model, to a potentially different space where the observation resides.

Ensemble methods are those that use many background states, each starting

from initial conditions that are slightly perturbed with respect to each other, rather than a single background state. This provides an easy way to estimate the background error covariance and also allows for it to evolve in time, better representing the "errors of the day." [20] Given an ensemble of M states denoted

$$\{\mathbf{x}^{b(1)}(t), \mathbf{x}^{b(2)}(t), \dots, \mathbf{x}^{b(M)}(t)\}, \quad (2.7)$$

the mean and ensemble perturbations about the mean are given by

$$\bar{\mathbf{x}}^b(t) = \frac{1}{M} \sum_{l=1}^M \mathbf{x}^{b(l)}(t) \quad (2.8)$$

and

$$\hat{\mathbf{X}}^b(t) = [\mathbf{x}^{b(1)}(t) - \bar{\mathbf{x}}^b(t) \quad \mathbf{x}^{b(2)}(t) - \bar{\mathbf{x}}^b(t) \quad \dots \quad \mathbf{x}^{b(M)}(t) - \bar{\mathbf{x}}^b(t)] \quad (2.9)$$

respectively. The background error covariance can be computed as

$$\mathbf{P}^b(t) = (M - 1)^{-1} \hat{\mathbf{X}}^b(t) \left(\hat{\mathbf{X}}^b(t) \right)^T. \quad (2.10)$$

A computationally efficient way of solving the cost function to determine the analysis ensemble is the Ensemble Transform Kalman Filter [22]. In this approach, the spread of the background ensemble about the mean is treated as a proper basis (although it is not, in fact, full rank). Arbitrary states can be expressed in this basis as $\mathbf{x} = \bar{\mathbf{x}}^b + \mathbf{w} \hat{\mathbf{X}}^b$ where $\mathbf{w} \in \text{span}(\hat{\mathbf{X}}^b)$. Defining the quantities $\bar{\mathbf{y}}^b = \mathbf{H} \bar{\mathbf{x}}^b$ and

$\hat{\mathbf{Y}}^b = \mathbf{H}\hat{\mathbf{X}}^b$, the cost function can be expressed as

$$\tilde{J}(\mathbf{w}) = \frac{1}{2}((M-1)\mathbf{w}^T\mathbf{w} + [\mathbf{y} - \bar{\mathbf{y}}^b - \hat{\mathbf{Y}}^b\mathbf{w}]^T (\mathbf{R}^o)^{-1} [\mathbf{y} - \bar{\mathbf{y}}^b - \hat{\mathbf{Y}}^b\mathbf{w}]). \quad (2.11)$$

The cost function is minimized to obtain the mean state

$$\bar{\mathbf{w}}^a = \tilde{\mathbf{P}}^a (\hat{\mathbf{Y}}^b)^T \mathbf{R}^{-1} (\mathbf{y}^o - \bar{\mathbf{y}}^b). \quad (2.12)$$

The members of the ensemble are given by the columns of

$$\hat{\mathbf{W}}^a = \left[(M-1) \left[(M-1)\mathbf{I} + (\hat{\mathbf{Y}}^b)^T \mathbf{R}^{-1} \hat{\mathbf{Y}}^b \right]^{-1} \right]^{1/2} \quad (2.13)$$

plus the mean $\bar{\mathbf{w}}^a$.

To obtain the analysis states, the mean is given by $\bar{\mathbf{x}}^a = \bar{\mathbf{x}}^b + \hat{\mathbf{X}}^b \bar{\mathbf{w}}^a$, and the members of the ensemble by $\mathbf{x}^{a(i)} = \bar{\mathbf{x}}^b + \hat{\mathbf{X}}^b \mathbf{w}^{a(i)}$. The analysis ensemble members are used as the initial conditions for the subsequent forecast for each background ensemble member respectively. Typically the forecast is made using a numerical model, however, in this case a data-derived phase space model of Sec. 2.2 will be used.

2.5 Nearest Neighbor Ensemble Transform Kalman Filter

The approach developed for this study combines the reconstructed phase space forecasting technique with the ETKF to produce improved forecasts. The first step

is to construct the model. This requires a long time series as discussed in 2.3. The model time series should not include any observations that will be used during the ETKF stage where forecasts of events will be made. In other words, all events to be forecast should be out of sample. The model construction proceeds as followed and is repeated a single time.

2.5.1 Model Construction

- (1) Given the model time series, compute the appropriate time delay τ and dimension m for the embedding.
- (2) Construct the time delay vectors as in Eq. 2.1 and the model trajectory \mathbf{D}^m .
- (3) To apply SSA, compute the eigenvalues and eigenvectors from the covariance matrix of the model trajectory as in Eq. 2.4.
- (4) Project the model trajectory onto the transform matrix of the k most important modes (Eq. 2.5).

2.5.2 Forecast Cycle

Once the model has been constructed, forecasts can be made using the ETKF and the model constructed from the time series data. The following sequence will be repeated to make forecasts after each forecast window of duration t_w . The initial background ensemble can be chosen from random points within the model trajectory or from an ensemble of nearest neighbors to an initial observation.

- (1) Beginning with an ensemble of prior estimates, compute the ensemble mean and spread.
- (2) Using the observations corresponding to time $\{t - (m-1)\tau, t - (m-2)\tau, \dots, t\}$, construct a time delay vector for the observation \mathbf{x}^o . To apply SSA, project the observation vector onto \mathbf{T}^m .
- (3) Apply the ETKF to the background ensemble and the observation to obtain the analysis ensemble.
- (4) Make the subsequent forecast using the data-derived reconstructed phase space model
 - (a) For each member of the analysis ensemble, locate the index of the analog among the model state vectors by minimizing the distance.

$$t^{NN(i)} = \min_{t \in [1, N']} |\mathbf{x}^a(i) - \mathbf{x}^m(t)| \quad (2.14)$$

- (b) Follow the evolution of the analog over the duration of the forecast window to obtain the background ensemble member $\mathbf{x}^{b(i)} = \mathbf{x}^m(t^{NN(i)} + t_w)$
- (5) Using the new background ensemble, proceed to step (2) to repeat the cycle for the next forecast window.

2.5.3 Reconstruction of Original Time Series

The forecasts are made in the model space consisting of time delay vectors. In the instances where singular spectrum analysis is used, the dimension of this space is reduced by projecting the time delay vectors onto a subset of eigenvectors representing the modes of variability. In order to compare forecasts made in this space to the observed time series a reconstruction must be performed.

- (1) If SSA has been applied to filter observational noise, a time delay vector corresponding to the forecast time is obtained by taking a superposition of the eigenvectors weighted by their time amplitude, i.e. the corresponding component of the forecast vector.
- (2) The forecasted value of the time series is the one with the latest time index, i.e. the m -th component of the time delay vector.

To evaluate how well the forecast performs it can be compared to a reference forecast. The simplest forecast to perform given a time series is a persistence forecast. Persistence assumes that the value of the variable will remain the same throughout the forecast window. Thus, the value at the end of the forecast window will be the same as the initial value at the start of the window. For both the nearest neighbor ETKF forecast and persistence, the root mean square error with respect to the true value will be computed. A skill score compares the error in each forecast and quantifies how well the ETKF forecast performed relative to the reference [35].

$$SS = 1 - \frac{RMSE_{\text{NN ETKF}}}{RMSE_{\text{persistence}}} \quad (2.15)$$

A skill score close to one means that the nearest neighbor ETKF forecast is performing better than persistence. A negative skill score means that persistence is outperforming the nearest neighbor ETKF forecast.

2.5.4 Localization

In common data assimilation scenarios, such as NWP, the model state is usually quite large covering many state variables and many grid points. The background error covariance matrix governs how each variable in the state vector is updated based on variations in the other state variables. Often the connection is obvious, however sometimes spurious correlations can occur within the background error covariance. This might lead to variables that are spatially far away having a disproportionate effect on one another. To counter this, a localization strategy is employed.

2.5.5 Data Density

The utility of this method is greatly determined by the quality of data in the model. It is essential to have sufficient observational data from which to construct a model. The first concern is that every portion of the attractor be visited by the model trajectory. The second is that every portion be visited frequently enough that sufficiently close nearest neighbor analogs can be located for each ensemble member.

2.6 Multi-channel Extension

If a partial state vector is available as an observation rather than simply a scalar time series, the techniques discussed above can still be applied [9]. Let the observed state vector \mathbf{x} be given by $(x_1(t), x_2(t), \dots, x_L(t))^T$, where L is the number of observed variables, or channels. A time delay vector can be constructed as follows:

$$\tilde{\mathbf{x}}(t) = (x_1(t), x_1(t+\tau), \dots, x_1(t+(m-1)\tau), \dots, x_L(t), x_L(t+\tau), \dots, x_L(t+(m-1)\tau))^T \quad (2.16)$$

such that the time-delay vector is $(m \times L)$ -dimensional looks like a concatenation of L m -dimensional time-delay vectors for each variable. The trajectory matrix and covariance matrix can be constructed as in Sec. 2.3. The eigenvalues and eigenvectors are also obtained in the same way. The eigenvectors, like the time-delay vectors, will appear as a piecewise configuration of L individual vectors. To reduce the dimension of the model, select the k eigenvectors with the largest eigenvalues that explain sufficiently much of the variability of the original data and project the trajectory matrix and observations onto this basis.

Chapter 3: Bred Vectors in the Reconstructed Phase Space

3.1 Introduction

A hallmark of nonlinear systems is that nearby trajectories tend to diverge exponentially as they evolve in time. This can lead to large errors in forecasts made based on initial conditions with relatively small errors. This is particularly evident in numerical weather prediction where the system has very complex and highly nonlinear behavior. “Errors of the day” arise from local instabilities that vary both spatially and temporally and are potentially fast growing. Ensemble forecasts are used to quantify the uncertainty due to the difficulty in capturing the fast growing and localized instabilities which causes the separation, or spread, among members of the ensemble to increase as the neighboring trajectories diverge. The technique of bred vectors was developed to maximize the potential spread of ensemble forecasts by predisposing members to grow in the presence of instabilities [14, 15].

Bred vectors are a computationally efficient way to pick up on directions of rapid growth within the model because they depend on numerically computing only two model runs - a control run and a perturbed run over a short breeding window. The direction of displacement between the perturbed trajectory and the control trajectory at the end of the breeding window is used to inform the choice of perturbed

initial condition for the next breeding cycle. After repeated cycles, bred vectors tend to pick up on the fastest direction of growth locally. It is natural to draw a comparison between the growth rate of bred vectors and Lyapunov exponents. Bred vectors converge to the local leading Lyapunov vector and their growth rate is a nonlinear extension of the leading local Lyapunov exponent [15].

Here we extend the technique of bred vectors, described below, to accommodate systems for which numerically integrating perturbed trajectories is not possible. The new technique, called Nearest-Neighbor Bred Vectors, relies on reconstruction of the phase space from a scalar time series, as described in Sec. 2.1. To evaluate the performance of the nearest-neighbors bred vectors technique in the phase space reconstructed from time series data, we turn to three simple, autonomous chaotic models – the Lorenz three variable model, the Chua oscillator, and the Rössler model. Each of these is often used as a test bed for nonlinear analysis techniques. These particular systems were chosen because they have two clearly defined regimes for which the signature of the transition between the two can be seen in the time series of a single model variable. We compare how well the original bred vectors can capture the transition between the two regimes in each system and test some simple prediction rules. We then compare the performance to nearest neighbor bred vectors in both the full phase space of all model variables and the reconstructed phase space of a single model variable.

3.1.1 Method of Bred Vectors

Given a set of model equations $\frac{d\mathbf{x}}{dt} = \mathbf{F}(\mathbf{x}(t))$, a control trajectory $\mathbf{x}_c(t)$ is obtained by integrating the model equations for a period of time. The breeding cycle starts by integrating a perturbed trajectory initiating from a point given by $\mathbf{x}_p(t_i) = \mathbf{x}_c(t_i) + \delta\mathbf{x}_i^0$, where $\delta\mathbf{x}_i^0$ is a displacement vector with a small magnitude equal to δ^0 . Both trajectories are evolved over an integration window of $n\Delta t$ where Δt is the time step. At the end of the integration window, the magnitude of the displacement vector between the perturbed trajectory and the control trajectory, given by $\delta\mathbf{x}_i^f = \|\mathbf{x}_c(t_{i+n}) - \mathbf{x}_p(t_{i+n})\|$ is compared to the initial displacement and the rate of exponential growth is computed as

$$g_i = \frac{1}{n\Delta t} \ln \left(\|\delta\mathbf{x}_i^f\| / \|\delta\mathbf{x}_i^0\| \right). \quad (3.1)$$

To start the next breeding cycle, the final perturbation vector is rescaled to the magnitude of the initial perturbation δ^0 such that the direction of the displacement final displacement from the previous cycle is preserved. The next breeding cycle initiates from a perturbed initial condition given by $\mathbf{x}_p(t_{i+n}) = \mathbf{x}_c(t_{i+n}) + \delta^0 \frac{\delta\mathbf{x}_i^f}{\|\delta\mathbf{x}_i^f\|}$. After repeated breeding cycles, the bred vectors will converge upon the local direction of greatest growth and their growth rates will approximate the local leading Lyapunov exponents [15]. This is a computationally efficient way to approximate similar information to that provided by Lyapunov exponents, but rather than computing a limit, only two integrations of the nonlinear model are required – one for

the control trajectory and one for the perturbed.

3.1.2 Nearest-Neighbor Bred Vectors

Many systems exist for which no numerical model is available. In this case we turn to a data-driven phase space model described in Ch. 2 and present an extension to the original breeding technique. Since this technique uses nearby points to serve as analogs for the perturbed initial conditions we call it nearest-neighbor breeding. Using a long historical time series of observations of a single variable, a phase space model is constructed as described in Sec. 2.1. The control trajectory is an out-of-sample time series of the same variable observed from the same system. To make forecasts of perturbations with respect to the control, analogs for the perturbed initial conditions are located among the model data points, the evolution of which is known. This analog is chosen such that the distance between the target perturbation and the analog is small, and the angle between the target perturbation and the nearest-neighbor perturbation is also small. This ensures that the direction of growth is maintained as much as possible during the breeding cycles.

From the time series data of a single variable, multivariate time delay coordinates are constructed according to Eq. 2.16. In this case, the control trajectory and model trajectory are each defined by discrete points. The control trajectory is given by \mathbf{x}_i^c and the perturbed trajectory is chosen from the model trajectory such that $\mathbf{x}_i^p = \mathbf{x}_j^m$.

The growth rate is computed as in Eq. 3.1 where $\delta\mathbf{x}_i^0 = \mathbf{x}_j^m - \mathbf{x}_i^c$ and

$\delta \mathbf{x}_i^f = \mathbf{x}_{j+n}^m - \mathbf{x}_{i+n}^c$ are the initial and final perturbations of the breeding cycle respectively. To select the perturbed trajectory of the next breeding cycle around the control starting from \mathbf{x}_{i+n}^c , we follow the spirit of the standard breeding in which the initial perturbation is given by $\delta \mathbf{x}_{i+n}^0 = \alpha \delta \mathbf{x}_i^f / \|\delta \mathbf{x}_i^f\|$; i.e., the final perturbation of the previous cycle is rescaled and bred as the new perturbation. Here the perturbation size α is constant for all breeding cycles. In the reconstructed phase space, however, the trajectory is defined by discrete points, and the rescaled position $\mathbf{x}_{i+n} + \alpha \delta \mathbf{x}_i^f / \|\delta \mathbf{x}_i^f\|$ may not be a trajectory point. We thus search and select the nearest point $\mathbf{x}_{j^*}^c$ and refer to the distance between these two points as the displacement distance.

Like the standard breeding, it involves two parameters, i.e. the window size n and the target perturbation size α . For successful applications of the nearest neighbor search, the density of the trajectory points must be high enough that, on average, the displacement distance in the nearest-neighbor search is small with respect to target initial perturbation size and the correlation between $\delta \mathbf{x}_i^f$ and $\delta \mathbf{x}_{i+n}^0$ is nearly 1.

3.1.3 Experimental Setup

The performance of nearest-neighbor bred vectors, particularly in the phase space reconstructed from time series data, is analyzed by performing the following three breeding experiments. Experiment (a) is the standard breeding in the 3-dimensional model phase space using the model equations as in Evans et al. [16].

Experiment (b) is the nearest-neighbor breeding applied to a discrete time series data in the model phase space, i.e., $\hat{\mathbf{x}}_i = \{x_1(t_i), x_2(t_i), \dots, x_n(t_i)\}$. This experiment reveals whether the nearest-neighbor breeding, without any knowledge about the model equations, gives comparable results to the standard breeding in the original model phase space. Finally, Experiment (c) is the nearest-neighbor breeding technique in the phase space reconstructed by the time-delayed embedding of a *single* time series i.e., $\tilde{\mathbf{x}} = \{x_j(t_i), x_j(t_{i+\tau}), \dots, x_j(t_{i+(m-1)\tau})\}$.

In all experiments, the growth rate g_i is computed using the breeding window size $n = 8$ with $\Delta t = 0.01$ and (targeted) perturbation size $\alpha = 0.10$. For each model, a control trajectory consists of approximately 10 000 orbits after an initial spin-up to make sure that the trajectory has reached the attractor. To ensure sufficient data density for the nearest-neighbor search in Experiments (b) and (c), nearest neighbors are selected from an out of sample model dataset that contains 800 000 data points in the original phase space.

3.2 Bred Vectors in the Reconstructed Phase Space of the Lorenz System

As a test bed for the techniques developed for this study, and here for the nearest-neighbors bred vectors, we use the Lorenz model to demonstrate the concepts of time series analysis, nonlinear dynamics, and data-driven data assimilation. The model was developed by Edward Lorenz in 1963 to serve as a idealized hydrodynamic model [36]. This set of first order ordinary differential equations is simple

in formulation,

$$\begin{aligned}\frac{dx}{dt} &= -\sigma x + \sigma y \\ \frac{dy}{dt} &= -xz + rx - y \\ \frac{dz}{dt} &= xy - bz,\end{aligned}\tag{3.2}$$

yet yields interesting nonlinear behavior.

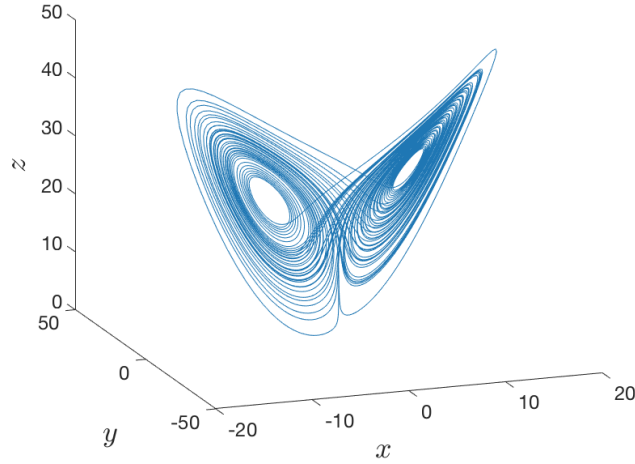


Figure 3.1: The attractor of the Lorenz 3 variable system in the natural 3 dimensional phase space of (x, y, z) .

The standard choice of parameters, $\sigma = 10$, $b = 8/3$, and $r = 28$, and integration with a time step of $\Delta t = 0.01$ produces the familiar butterfly attractor seen in Fig. 3.1. While the system has only three model dimensions, the strange attractor onto which the phase space collapses has a fractional dimension between 2 and 3. Three fixed points organize the topology of the attractor. Two foci, both unstable spiral points, give rise to the two lobes about which the trajectory orbits. A hyper-

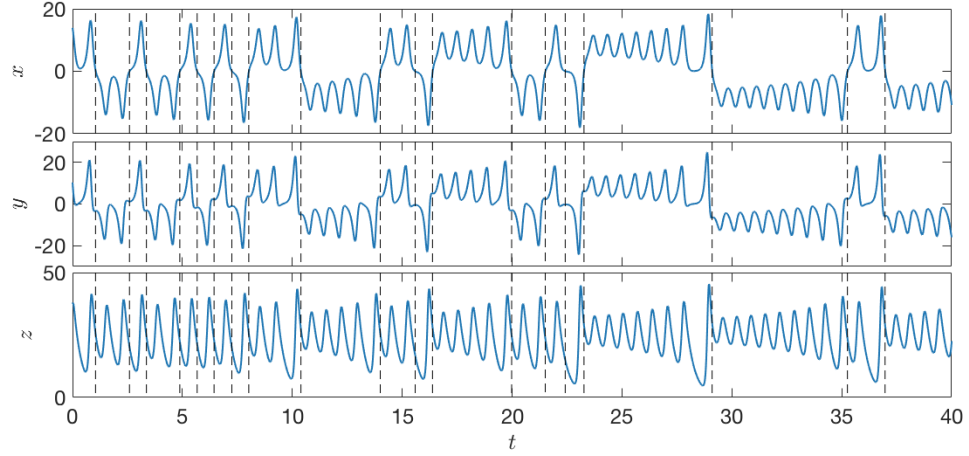


Figure 3.2: The time series data for the variables $x(t)$, $y(t)$, and $z(t)$ for the Lorenz system. The signature of regime change can be seen in the $x(t)$ variable when a sign change occurs.

bolic fixed point located at the origin governs the chaotic transition between two regimes. The two regimes are defined by the two lobes of attractor. The signature of the transition can be seen in positive or negative values of the x coordinate. After a number of orbits about the fixed point in one of the lobes, the trajectory crosses the $x = 0$ axis to orbit the other lobe. The behavior is nonperiodic and the number of orbits varies irregularly revealing the chaotic nature of the system. Though the behavior of the system is chaotic, it is predictable, making it a unique nonlinear system, ideal for testing new predictive techniques [16, 23, 37, 38].

3.2.1 Reconstruction of the Lorenz Attractor

To reconstruct the Lorenz attractor from the time series $x(t)$, time delay vectors are constructed as described in Sec. 2.1. The information shared by the first

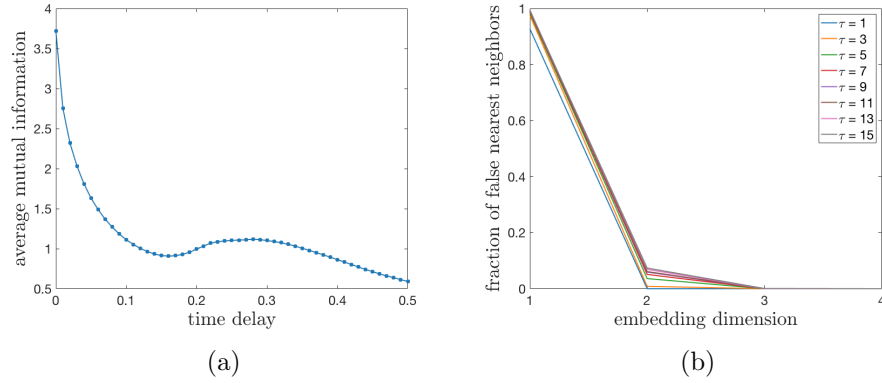


Figure 3.3: 3.3(a) The mutual information function as a function of time delay for the Lorenz $x(t)$ variable. 3.3(b) The percentage of false nearest neighbors for attractors reconstructed from the $x(t)$ variable of the Lorenz system for various time delays as a function of embedding dimension. For a range of time delays, the percentage of false nearest neighbors becomes negligible for more than three dimensions.

variable of the Lorenz system, $x(t)$, and its time-lagged counterpart, $x(t + \tau\Delta t)$, reaches a first minimum after a time delay of $t = 0.14$ or $\tau = 14$ time steps of $\Delta t = 0.01$, as seen in Fig. 3.3(a). At the first minimum, $x(t)$ and $x(t + \tau\Delta t)$ share relatively little information with one another, thus introducing maximal information to the vector. The time required for the information shared by the variable to decay to e^{-1} of its original content at $t = 0.07$. Often the embedding is not particularly sensitive to the value of the time delay chosen and choices of τ within this range will yield a proper embedding. For the reconstruction of the Lorenz attractor from the $x(t)$ variable a time delay of $\tau = 7$ will be used.

The minimum embedding dimension of the time delay vectors is chosen such that the attractor is unfolded and the phase space trajectories do not intersect. The method of false nearest neighbors described in Sec. 2.2 quantifies how neighborhoods of points on the attractor change as the embedding dimension is increased. If phase space trajectories intersect, points that should not be nearby will be spuriously close.

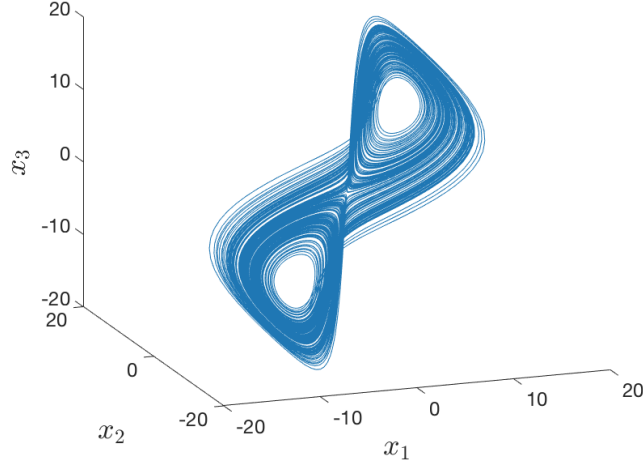


Figure 3.4: The reconstructed attractor in the embedded phase space formed by time delay vectors with dimension $m = 3$ and time delay $\tau = 7$.

As a sufficient number of dimensions is reached, points occupying neighborhoods should stabilize and remain unchanged. For all of the selected time delays, the percentage of false nearest neighbors shown in Fig. 3.3(b) tends to zero after $m = 3$, indicating a minimum of three dimensions are required for a proper embedding of the Lorenz attractor.

The reconstructed Lorenz attractor is depicted in Fig. 3.4. While the form of the attractor appears different than the full phase space attractor in Fig. 3.1, the major features are retained. Because the topological properties of the attractor are preserved by the skeleton of the attractor in the reconstructed phase space, the regime change between orbiting one foci of the butterfly to the other governed by the hyperbolic fixed point can be seen in each.

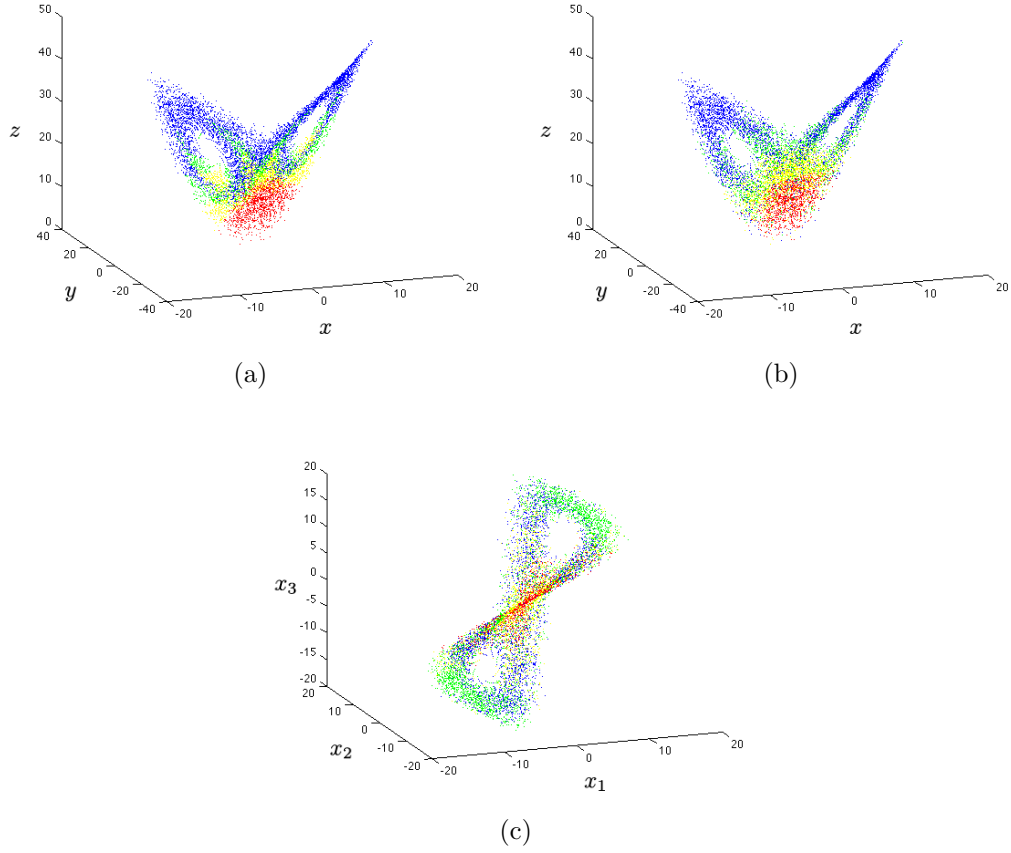
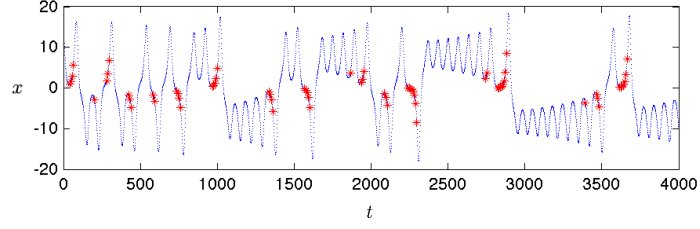


Figure 3.5: Growth rates of bred vectors in the Lorenz system using three different methods: 3.5(a) standard breeding using the ordinary differential equations in the phase space (x, y, z) ; 3.5(b) nearest-neighbor breeding in the phase space (x, y, z) ; and 3.5(c) nearest-neighbor breeding in the reconstructed phase space (x_1, x_2, x_3) . The colored points correspond to negative (blue), low (green), medium (yellow) and high (red) growth; see text for the value of the thresholds.

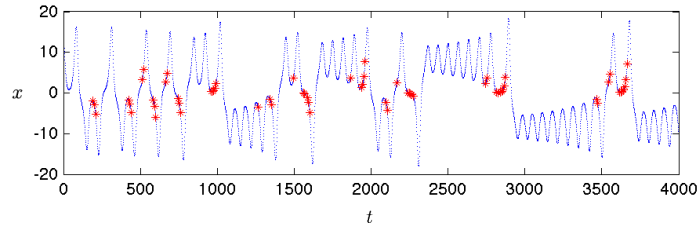
3.2.2 Bred Vectors Results

The average displacement distance and vector correlation between the standard and the nearest neighbor breeding are 0.17 and 0.97 in Experiment (b); they are 0.12 and 0.98 in Experiment (c).

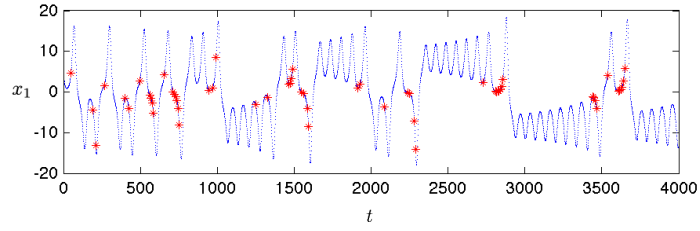
Figure 3.5 shows the points along the respective control trajectories for which bred vector growth rates were computed in each of the three experiments. The



(a)



(b)



(c)

Figure 3.6: The first coordinate of phase space as a function of time, with red stars indicating the points with high growth rate ($g_i \geq 6.4$) bred vectors; (a)-(c) are the same as in Fig. 3.5.

magnitudes of the growth rate are represented by colors using the same empirical thresholds as in Evans et al. [16]: negative growth points ($g_i < 0$) in blue, low growth points ($0 \leq g_i < 3.2$) in green, medium growth points ($3.2 \leq g_i < 6.4$) in yellow, and high growth points ($g_i \geq 6.4$) in red. As shown in Evans et al. [16] for the standard breeding, all experiments show high growth at points concentrated in the regime transition region, while the regions with different growth rates are well separated. The nearest-neighbor breeding, both in the original (Fig. 3.5b) and in the reconstructed (Fig. 3.5c) phase spaces, successfully captures the features found in the standard breeding (Fig. 3.5a), although the separation between the different growth rates is less sharp. Figure 3.6 shows the time series of the first phase space coordinate (x) for the first 500 breeding cycles for each of the three experiments. Note that, by construction, the first coordinate x in phase space coincides with the first coordinate x_1 in the embedded space.

3.2.3 Predicting Regime Changes

As pointed out by Evans et al. [16] and apparent in Fig. 3.6, high bred vector growth rate, marked in red, is a very good predictor of regime change in the standard breeding in the Lorenz model. To test the predictive capabilities of the bred vector growth rates, we make binary (Yes-No) forecasts based on the prediction rules developed by Evans et al. [16]:

- **Rule 1:** The appearance of a red star indicating a growth rate in excess of 6.4 signals a regime change will occur at the end of the

current orbit (Yes). The absence of a red star means the current regime will continue (No).

- **Rule 2:** The length of the new regime is proportional to the number of red stars. For example, the presence of five or more stars in the old regime, implies that the new regime will last four orbits or more (Yes). Fewer than five stars means the new regime will last fewer than four orbits (No).

We modify these rules slightly for the case of nearest-neighbor breeding by excluding red stars that occur on orbits with extrema whose absolute value is greater than 1, to reduce the false alarms in the prediction of regime change.

Table 3.1 is the contingency table [35] of the forecast/observed event pairs applying the first prediction rule to the bred vectors from the three experiments. Individual forecasts of the rule (FCST) are made for each orbit of the trajectory. The observed events (OBS) are based on the occurrence or non-occurrence of the regime change following the completion of the orbit. Corresponding accuracy measures [35] are shown in Table 3.2.

Examining the hit rate (HR), threat score (TS), and false alarm rate (FAR), it is apparent that the three experiments succeed in predicting the regime change with similar accuracy. The HRs and TSs for the three methods are close, varying from 82 to 87, and 72 to 76 %, respectively. The FARs are about 6% for the standard breeding but increase to 11 and 13% when nearest-neighbor breeding is used in the original model phase space and in the reconstructed phase space, respectively.

Table 3.1: Contingency tables based on the rule that regime change will occur in the orbit following the appearance of high growth rate bred vectors using three different methods. In (b) and (c) using the nearest-neighbor breeding, high growth rate points in orbits with absolute values of extrema above 1 are excluded. OBS and FCST stand for observed and forecast, respectively; (a)-(c) are the same as in Fig. 3.5.

			OBS		
			Yes	No	Total
(a)	FCST	Yes	374	38	412
		No	80	573	653
		Total	454	611	1065
(b)	FCST	Yes	396	67	463
		No	58	544	602
		Total	454	611	1065
(c)	FCST	Yes	383	77	460
		No	71	534	605
		Total	454	611	1065

We note that, in addition to large bred vector growth rate, two other methods have been also proposed to predict regime changes in the Lorenz three variable system. In his original paper (Lorenz, 1963), Lorenz pointed out that regime changes were associated with large values of the variable z . Yadav et al. (2005) showed that large absolute magnitudes of the x variable are also a good predictor. We have implemented the method used in [38] for the time series $x(t)$ and got equally good results. However, the main objective of this paper is to determine whether bred vectors can predict stability from a single time series data, i.e., in the reconstructed phase space. The reasons for the choice of bred vectors as a predictor of fast growth in the dynamical system are two-fold. First, unlike the size of a particular variable,

Table 3.2: Measures of forecast accuracy in terms of the Hit Rate (HR), Threat Score (TS), and False Alarm Rate (FAR); (a)-(c) are the same as in Fig. 3.5. The final row shows the values when the threshold of $x(t_i)$ rule is used.

	HR (%)	TS (%)	FAR (%)
(a)	82.4	76.0	5.2
(b)	87.2	76.0	11.0
(c)	84.4	72.1	12.6

breeding can be tested in any dynamical model. Second, bred vector perturbations and their growth have a clear physical meaning in that they detect instabilities [39] and are akin to the leading local Lyapunov vector and their corresponding growth [40]. Thus while predictions based on threshold values of a single variable work well for the Lorenz model, bred vector growth rate may be suitable for making predictions in a broad range of dynamical systems. The bred vectors, by their ability to capture nonlinearity in instability growth, can characterize instability in dynamical systems in a more general manner than Lyapunov vectors and provide a way to obtain regime change in cases where the latter are computed from the time series data [41].

For the binary forecasts, Evans et al. (2004) noted that the next regime tended to be longer-lasting when the when more red stars appeared in the previous regime. As can be seen in Fig. 3.7, this rule holds true for the nearest neighbor breeding technique in the full phase space. However, it is less convincing for high growth rate bred vectors in the reconstructed phase space.

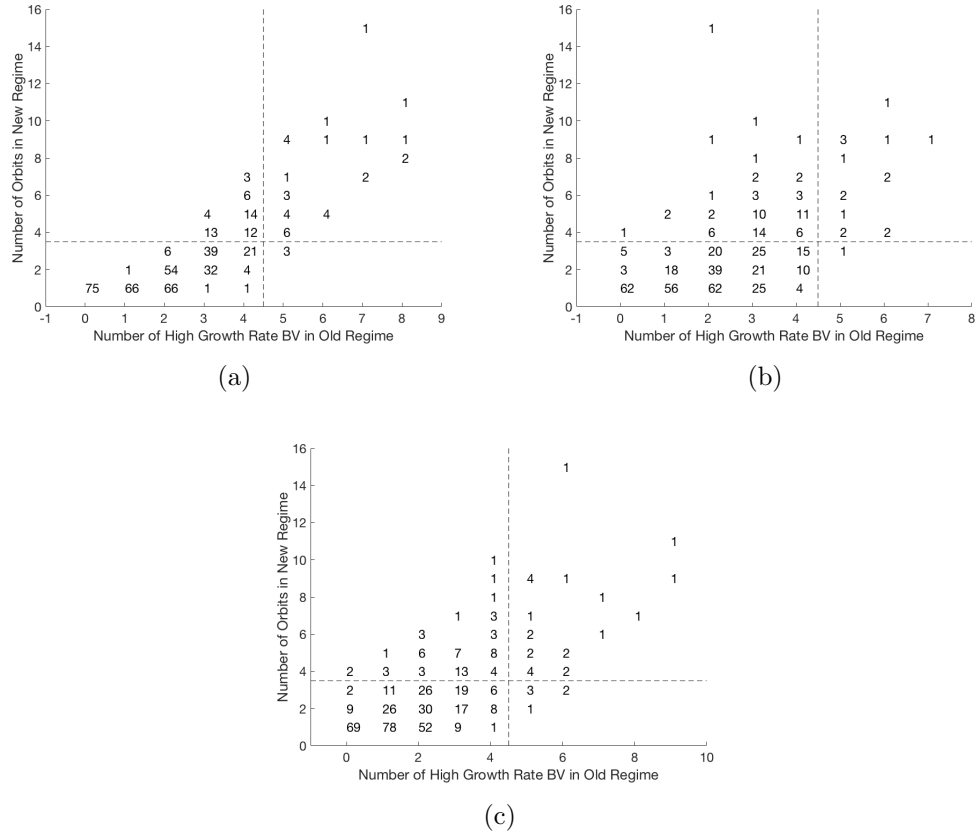


Figure 3.7: The observed number of cycles vs the number of high growth rate bred vectors observed in previous regime. The numbers indicate the frequency with which a given pair was observed. The line correspond to the values in prediction rule 2. (a) Bred vectors are computed using the ODEs. (b) NN bred vectors are computed using the full phase space model. (c) NN bred vectors are computed using the reconstructed phase space model.

3.3 Bred Vectors in the Reconstructed Phase Space of the Chua Oscillator

The Chua oscillator is one of the first electronic circuits constructed to display chaotic behavior [42]. This system produces a double scroll attractor, similar to the Lorenz attractor, with two lobes that the trajectory orbits and a transition between the two regimes, as seen in Fig. 3.8. However, unlike the Lorenz system,

the equations that model the behavior do not contain an explicit nonlinear term.

The system is described by the following equations:

$$\begin{aligned}\frac{dx}{dt} &= \alpha (y - x - h(x)) \\ \frac{dy}{dt} &= x - y + z \\ \frac{dz}{dt} &= \beta y,\end{aligned}\tag{3.3}$$

where $h(x) = bx + \frac{1}{2} (a - b) (|x + 1| - |x - 1|)$. The attractor in Fig. 3.8 is produced using the values of $\alpha = 11.6$, $\beta = 18.432$, $a = -1.4554$, and $b = -0.7853$ for the coefficients.

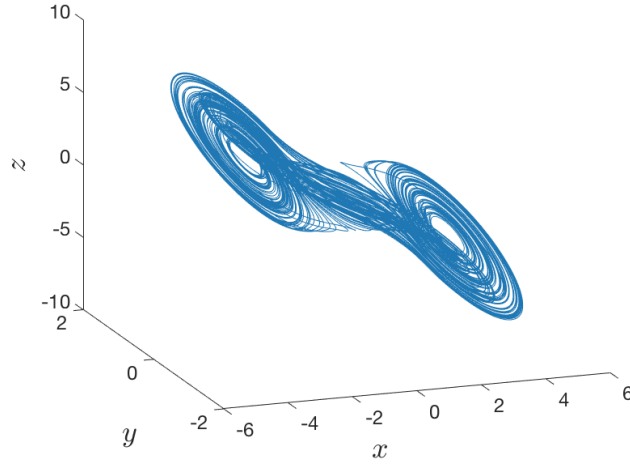


Figure 3.8: The attractor of the Chua 3 variable system in the natural 3 dimensional phase space of (x, y, z) .

Despite the lack of an explicitly nonlinear term in the governing equations, the system has rigorously been shown to exhibit chaotic behavior [43]. In fact, the system was constructed to mimic the type of double scroll behavior seen in the Lorenz system. The Chua attractor has been used in numerous studies as a model

testbed for nonlinear techniques.

The three fixed points of the system divide the domain into three regions where the behavior can be analyzed. Like the Lorenz system, there are two fixed points with a stable and unstable manifold that results in a trajectory that orbits each. The one-dimensional stable manifold draws trajectories into the plane spanned by the unstable manifold, where trajectories tend to spiral away from the fixed point. These orbits can be seen in the time series of both the x and y coordinates, with one region defined by positive x values and negative y values, and the other by negative x and positive y coordinates. Also similar to the Lorenz system, the Chua circuit attractor possesses a fixed point at the origin that controls the transition between the two regimes. However, this fixed point is not hyperbolic.

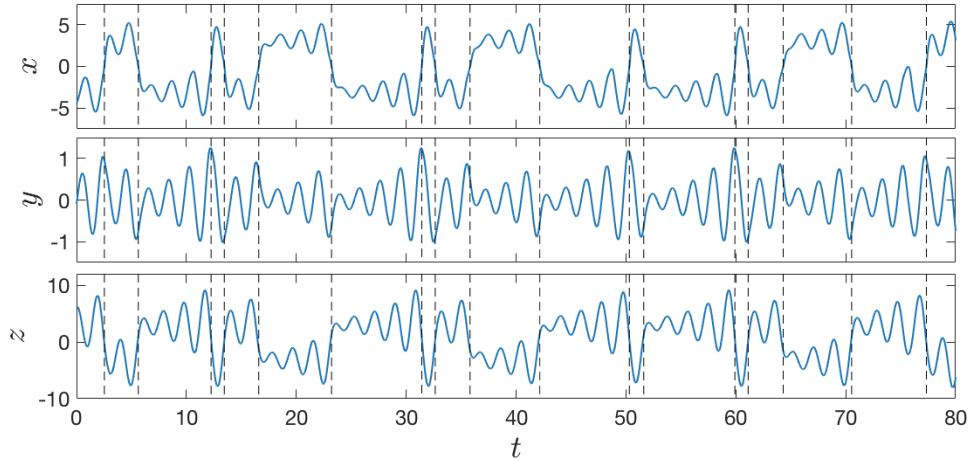


Figure 3.9: The time series data for the variables $x(t)$, $y(t)$, and $z(t)$ for the Chua system. The signature of regime change can be seen in the $x(t)$ variable when a sign change occurs.

The signature for the regime change can be seen in the value of the x coordinate

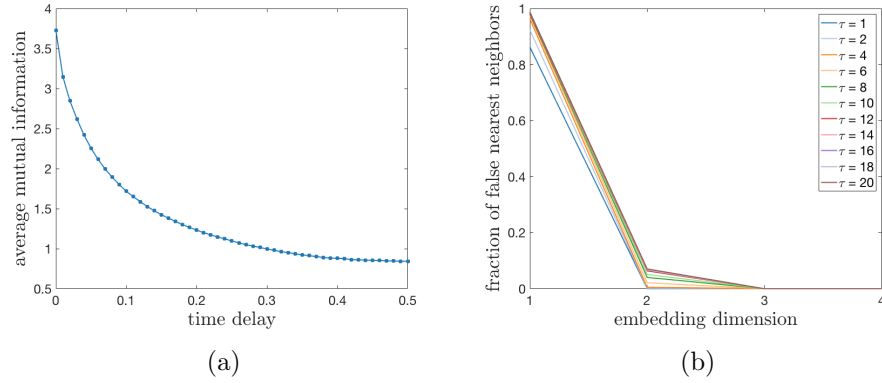


Figure 3.10: 3.10(a) The mutual information function as a function of time delay for the Chua $x(t)$ variable. 3.10(b) The percentage of false nearest neighbors for attractors reconstructed from the $x(t)$ variable of the Chua system for various time delays as a function of embedding dimension. For a range of time delays, the percentage of false nearest neighbors becomes negligible for more than three dimensions.

in Fig. 3.9. As it crosses from positive to negative and vice versa, the trajectory transitions from orbiting the fixed point in the positive region to that in the negative region.

3.3.1 Reconstruction of the Chua Attractor

Time delay vectors built from the first variable of the Chua system are used to reconstruct the attractor. The mutual information function in Eq. 2.2 computed for a series of time lags is shown in Fig. 3.10(a). Unlike in that of the Lorenz $x(t)$ variable, the average mutual information of the Chua $x(t)$ variable displays no local minima. However, the e -folding time of the mutual information of about $\tau = 16$ can be used to reconstruct the attractor.

The fraction of false nearest neighbors, shown in Fig. 3.10(b), falls off to negligible values after three dimensions. The attractor reconstructed from time delay vectors with delay of 16 time steps between each component and three dimensions

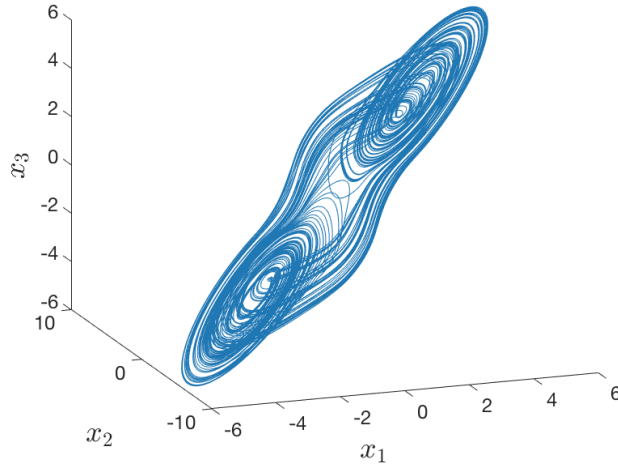


Figure 3.11: The reconstructed Chua attractor in the embedded phase space of time delay vectors with $m = 3$ and $\tau = 16$.

is shown in Fig. 3.11. The double scroll of the Chua attractor in Fig. 3.8 can still be seen. The transition region about the origin is also preserved.

3.3.2 Bred Vectors Results

High growth rate bred vectors in the Lorenz system tended concentrate about the hyperbolic fixed point at the origin where the instability gave rise to regime change. In the Chua system, high growth rate bred vectors also concentrate near the fixed point at the origin; however, their appearance is not as clear a precursor to regime change as it was for the Lorenz system.

The distribution of bred vector growth rates are very similar for the standard bred vectors and those computed using the nearest neighbor bred vector method in the full (x, y, z) phase space, as can be seen in Fig. 3.12. Here, the points with negative growth rates are blue, low growth rate points ($0 \leq g_i < 2.85$) in green, medium growth points ($2.85 \leq g_i < 5.7$) are yellow, and high growth rate points

($g_i \geq 5.7$) are red. In the standard and full phase space nearest neighbor bred vectors, the high growth rate points are mostly in the transition regions, however, there are some instances of high growth rate bred vectors in the scroll portions of the attractors.

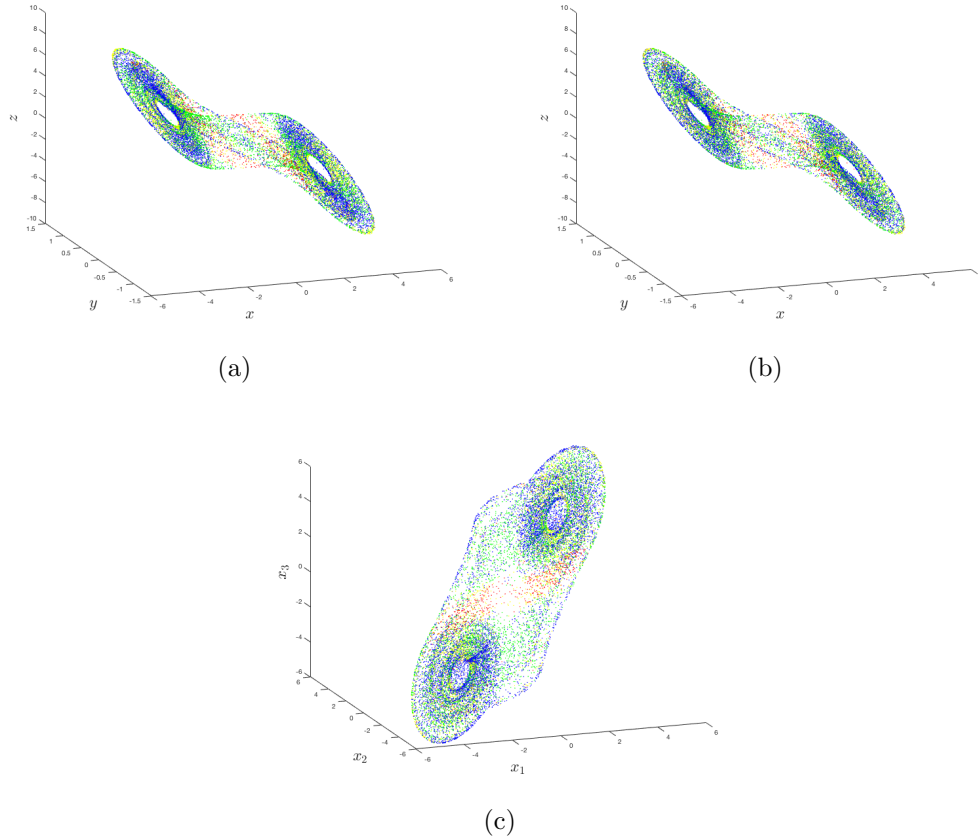
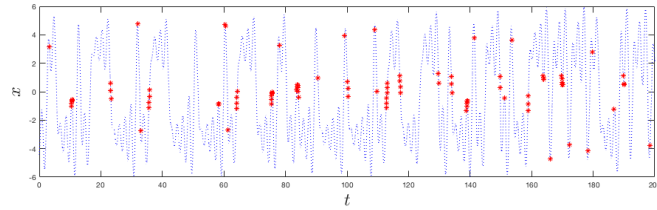


Figure 3.12: Growth rates of bred vectors in the Chua system using three different methods: 3.12(a) standard breeding using the ordinary differential equations in the phase space (x, y, z) ; 3.12(b) nearest-neighbor breeding in the phase space (x, y, z) ; and 3.12(c) nearest-neighbor breeding in the reconstructed phase space (x_1, x_2, x_3) . The colored points correspond to negative (blue), low (green), medium (yellow) and high (red) growth; see text for the value of the thresholds.

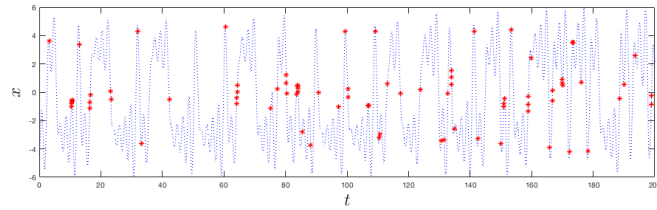
The same threshold values are used for bred vectors in the reconstructed phase space. Here the separation between the growth rates and the occurrence of high growth rate bred vectors is even more clear, which contrast with the results seen in

the same scenario using the Lorenz system where the clearest divisions were seen in the standard breeding technique using the ODEs.

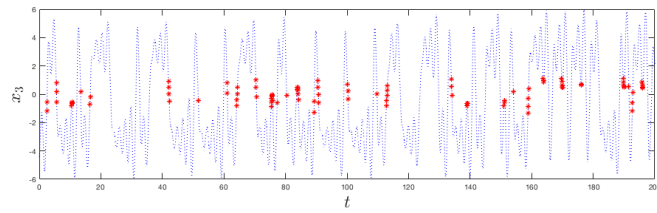
This result can also be seen in the time series in Fig. 3.13. In the time series from experiments (a) and (b), occasionally high growth rate bred vectors can be seen near the extreme points of the orbits. However, these points do not always occur in the final orbit before the transition. In experiment (c) the high growth rate bred vectors rarely occur close to the extrema of the orbits and are confined almost exclusively to the transition region.



(a)



(b)



(c)

Figure 3.13: The first coordinate of phase space as a function of time, with red stars indicating the points with high growth rate ($g_i \geq 6.4$) bred vectors; (a)-(c) are the same as in Fig. 3.12(a).

3.3.3 Predicting Regime Change

While the high growth rate bred vectors are not as strong a predictor of regime change as in the case of the Lorenz system, they still provide good skill for forecasting regime change. Table 3.3 shows the contingency table for 967 orbits for prediction Rule 1, that the appearance of high growth rate bred vectors indicate that the orbit will be the final one in the regime.

Table 3.4 contains the skill scores for forecast rule 1. Given the threshold values for bred vector growth rates defined above, all three methods produce forecasts with similar skill. The nearest neighbor bred vectors in the full model phase space had the highest percentage of false alarm rates indicating a high growth rate bred vector did not appear despite impending regime change.

The method of nearest neighbor bred vectors in the reconstructed phase space had both the highest percentage of correct forecasts and the lowest false alarm rate. In fact, with a lower threshold, the embedded bred vectors have the potential to predict an even higher percentage of regime changes correctly, with $PC = 81.6$, $HR = 75.5$, $TS = 63.4$, $FAR = 14.0$. For this system, predictions made using a time series of a single variable are even more skillful than methods where full knowledge of the equations of motion or phase space are utilized.

Table 3.3: Contingency tables based on the rule that regime change will occur in the orbit following the appearance of high growth rate bred vectors using three different methods. In (b) and (c) using the nearest-neighbor breeding, high growth rate points in orbits with absolute values of extrema above 1 are excluded. OBS and FCST stand for observed and forecast, respectively; (a)-(c) are the same as in Fig. 3.12.

			OBS		
			Yes	No	Total
(a)	FCST	Yes	242	69	311
		No	166	490	656
		Total	408	559	967
(b)	FCST	Yes	247	113	360
		No	161	446	607
		Total	408	559	967
(c)	FCST	Yes	308	78	386
		No	100	481	581
		Total	408	559	967

3.4 Bred Vectors in the Reconstructed Phase Space of the Rössler System

Another well known and well studied simple nonlinear system was developed by Otto Rössler in the late 1970s [44, 45]. The Rössler system also has been used as a nonlinear testbed for studies [46]. Constructed to be an even simpler flow than that of the Lorenz system, the attractor of the Rössler system has only one scroll.

Table 3.4: Measures of forecast accuracy in terms of the Hit Rate (HR), Threat Score (TS), and False Alarm Rate (FAR); (a)-(c) are the same as in Fig. 3.12.

	PC (%)	HR (%)	TS (%)	FAR (%)
(a)	75.7	59.3	50.7	12.3
(b)	71.7	60.5	47.4	20.2
(c)	77.0	57.4	51.3	8.6

The equations have a single nonlinear term and are given by:

$$\begin{aligned}
 \frac{dx}{dt} &= -y - z \\
 \frac{dy}{dt} &= z + ay \\
 \frac{dz}{dt} &= b + z(x - c).
 \end{aligned} \tag{3.4}$$

The shape of the attractor varies based on the choice of parameters. Setting $a = 0.55$, $b = 2$, and $c = 4$ yields an attractor with a screw shape seen in Fig. 3.14. The system has two fixed points, located at

$$\left(\frac{c}{2} \pm \sqrt{c^2 - 4a}, -\frac{c}{2a} \mp \frac{1}{2a}\sqrt{c^2 - 4a}, \frac{c}{2a} \pm \frac{1}{2a}\sqrt{c^2 - 4a} \right). \tag{3.5}$$

The regime change between orbiting within in the xy -plane to approaching the unstable fixed point outside of the plane is chaotic. The signature for regime change can be seen in the x and z time series, however it is most clear in the $x(t)$ data. The regime change can be seen when the value of $x(t)$ crosses from positive to negative, similar to the Lorenz system. This coincides with transitioning from orbiting close to the $z = 0$ plane, where values of $z(t)$ tend to be close to zero, to orbiting the fixed

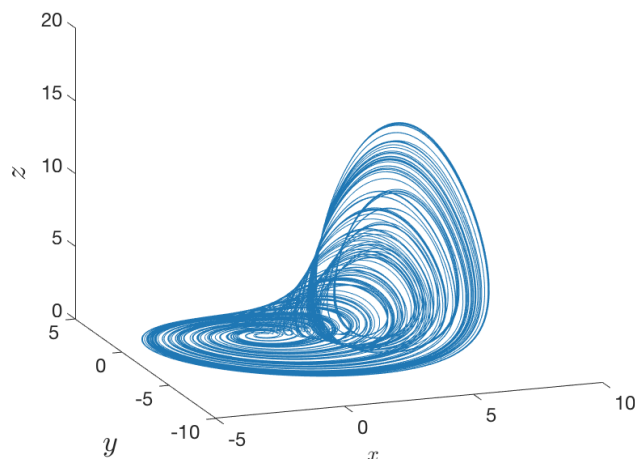


Figure 3.14: The attractor of the Rössler 3 variable system in the natural 3 dimensional phase space of (x, y, z) .

point outside of that plane in the screw-like portion of the attractor. One obvious difference between the regimes in this system, and those in the Lorenz and Chua systems, is that in the negative x regime, the system always completes a single orbit before transitioning to the positive x regime, as can be seen in Fig. [3.15](#)

3.4.1 Reconstruction of the Rössler Attractor

The Rössler attractor has three state variables and the attractor has a dimension between two and three. The fraction of false nearest neighbors falls off after three embedding dimensions for a variety of time delays chosen, as can be seen in Fig. [3.16\(b\)](#). The average mutual information function for the time series of the $x(t)$ variable, in Fig. [3.16\(a\)](#), falls off exponentially. As for the Chua system $x(t)$ variable, there are a range of appropriate time delays for the reconstruction ranging from a single timestep to approximately the e -folding time of the mutual information function. As the time delay increases, the orbits of the two regimes become more

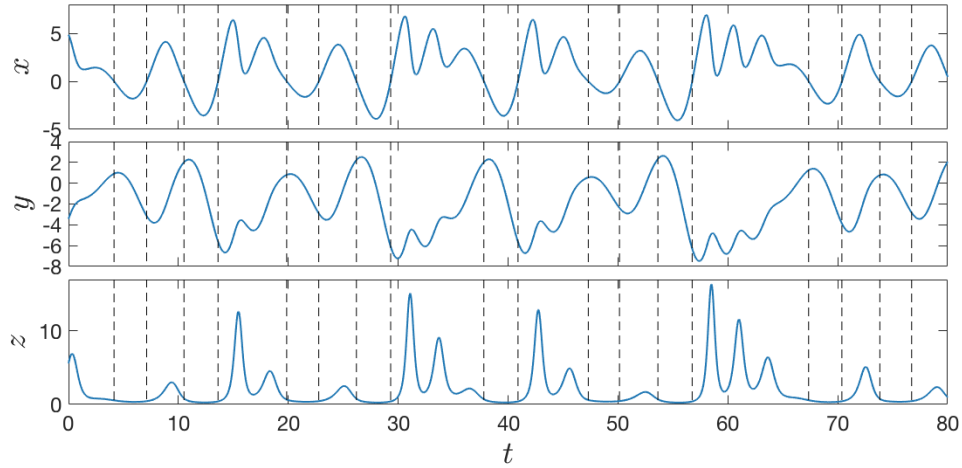


Figure 3.15: The time series data for the variables $x(t)$, $y(t)$, and $z(t)$ for the Rössler system. The signature of regime change can be seen in the $x(t)$ variable when a sign change occurs.

clearly defined. Here, a time delay of 34 time steps is chosen and the two regimes are visible in Fig. 3.17.

3.4.2 Bred Vectors Results

A striking result of the distribution of bred vector growth rates in the Lorenz system, using the ODEs to compute the evolution of the perturbed trajectory, is the clear separation between the regions of various growth rates. These regions became less clearly defined in the case of the nearest neighbor bred vectors in the reconstructed phase space. The same feature is evident for the Rössler system. Low growth rate bred vectors (indicated by green points) are found almost entirely in the $z = 0$ plane in Fig. 3.18(a). Yellow and red points representing medium and high growth rate bred vectors respectively are found just after transition from orbiting

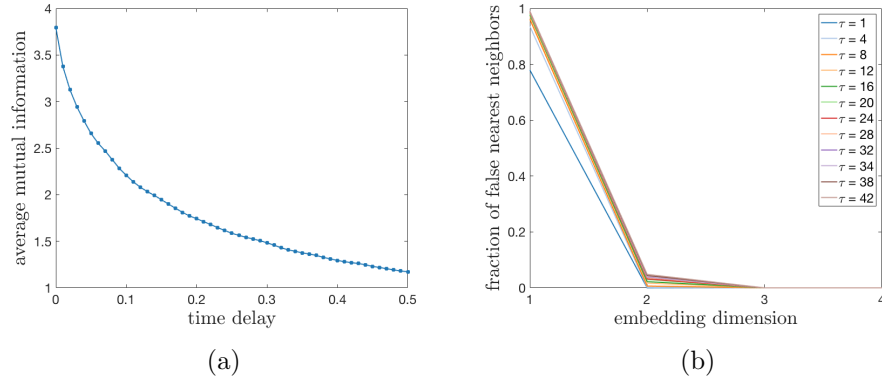


Figure 3.16: 3.16(a) The mutual information function as a function of time delay for the Rössler $x(t)$ variable. 3.16(b) The percentage of false nearest neighbors for attractors reconstructed from the $x(t)$ variable of the Rössler system for various time delays as a function of embedding dimension. For a range of time delays, the percentage of false nearest neighbors becomes negligible for more than three dimensions.

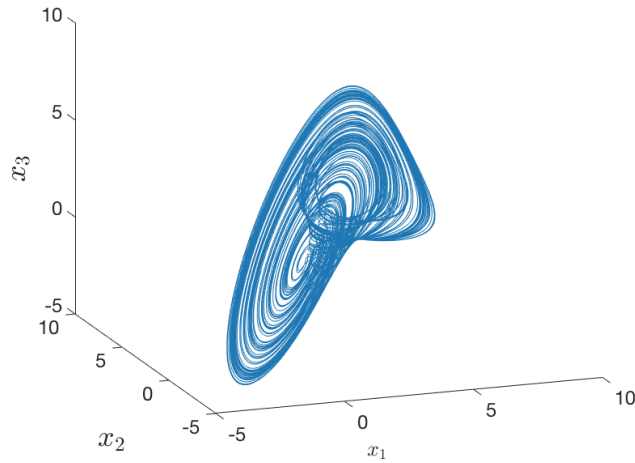


Figure 3.17: The reconstructed attractor in the embedded phase space formed by time delay vectors with dimension $m = 3$ and time delay $\tau = 34$.

the xy -plane regime to the screw portion of the attractor when the sign of the x variable changes. A similar pattern is seen in Fig. 3.18(b) for nearest neighbor bred vectors computed using the full (x, y, z) phase space to determine the evolution.

These features are not well defined for nearest neighbor bred vectors computed in the reconstructed phase space seen in Fig. 3.18(c). Regions of negative, low, and medium growth rate bred vectors occur throughout the attractor. High growth rate bred vectors, however, do tend to be confined to a single region and tend to precede the regime change. This can be seen more clearly in Fig. 3.18(c). Similarly, the fact that high growth rate bred vectors come immediately following the regime change from the negative to positive x regimes in the other two cases can be seen in Figs. 3.19(a) and 3.19(b). In these experiments, no high growth rate bred vectors appear in the negative x regime, therefore the transition from the negative to positive regime can not be predicted by bred vector growth rate. In the embedded phase space experiment, a high growth rate bred vector often appears in the negative regime.

3.4.3 Predicting Regime Change

Because the high growth rate bred vectors tend to occur immediately following a regime change in the cases where knowledge of the full phase space is used, they are poor predictors of regime change. Contingency tables for each of the three experiments are given in Table 3.5. Threshold values of 0.9 and 1.8 were chosen for these experiments, but no good choice of threshold for any of the three scenarios.

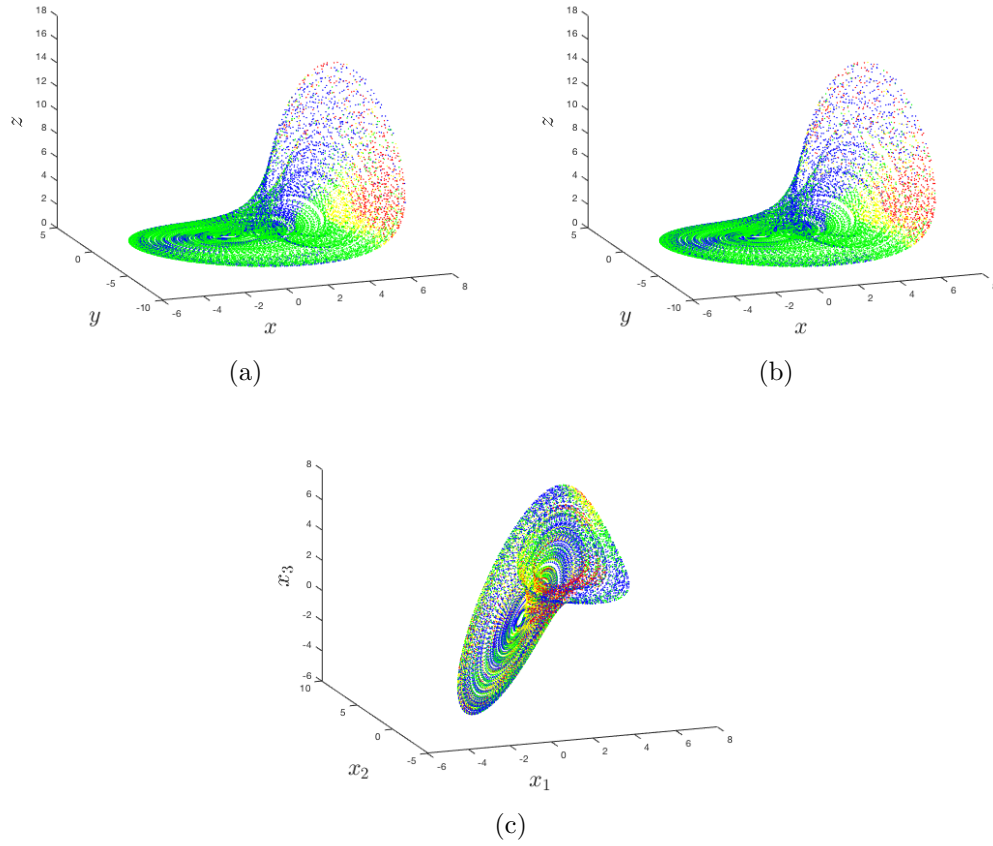
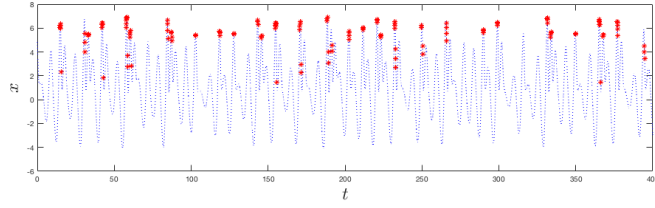


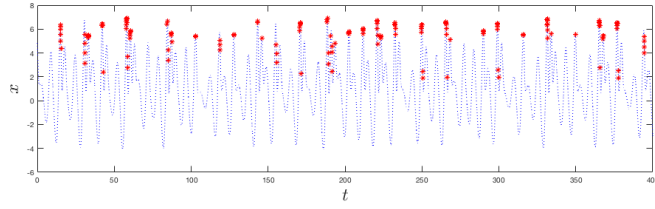
Figure 3.18: Growth rates of bred vectors in the Lorenz system using three different methods: 3.18(a) standard breeding using the ordinary differential equations in the phase space (x, y, z) ; 3.18(b) nearest-neighbor breeding in the phase space (x, y, z) ; and 3.18(c) nearest-neighbor breeding in the reconstructed phase space (x_1, x_2, x_3) . The colored points correspond to negative (blue), low (green), medium (yellow) and high (red) growth; see text for the value of the thresholds.

These values maximize the predictive skill in the case of the nearest neighbor bred vectors in the reconstructed phase space. While the ODE and nearest neighbor experiments with the full phase space show very little, if any, predictive skill, as indicated by the very small percentage of correct forecasts in Table 3.6, there is moderate predictive skill in the embedded phase space experiment.

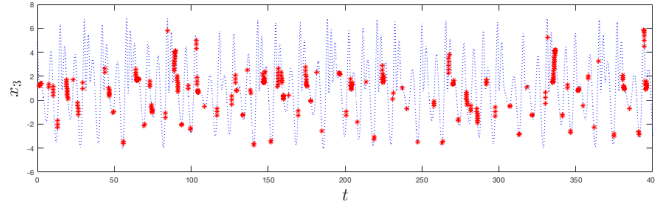
In the full phase space experiments, a high growth rate bred vector appears following every transition to the positive regime, however often this regime persists



(a)



(b)



(c)

Figure 3.19: The first coordinate of phase space as a function of time, with red stars indicating the points with high growth rate ($g_i \geq 6.4$) bred vectors; (a)-(c) are the same as in Fig. 3.18(a).

for more than one orbit, thus the high false alarm rate. There are no high growth rate bred vectors in the negative x regime. High growth rate bred vectors cannot predict the regime change from negative to positive resulting in an extremely low hit rate.

The predictions made using high growth rate bred vectors from nearest neighbor breeding in the reconstructed phase space are more successful. Regime changes from negative to positive and positive to negative values of x can both be forecasted. However, the skill is moderate compared to similar predictions made in the Lorenz and Chua systems.

Table 3.5: Contingency tables based on the rule that regime change will occur in the orbit following the appearance of high growth rate bred vectors using three different methods. In (b) and (c) using the nearest-neighbor breeding, high growth rate points in orbits with absolute values of extrema above 1 are excluded. OBS and FCST stand for observed and forecast, respectively; (a)-(c) are the same as in Fig. 3.18.

			OBS		
			Yes	No	Total
(a)	FCST	Yes	7	267	274
		No	725	31	756
		Total	732	298	1030
(b)	FCST	Yes	15	276	291
		No	717	22	739
		Total	732	298	1030
(c)	FCST	Yes	515	136	651
		No	217	162	379
		Total	732	298	1030

3.5 Summary and Conclusions

Bred vectors are useful because they provide an efficient way to detect instabilities. Previously, bred vectors have been shown to be useful in predicting regime change for the Lorenz system [16]. Here we have extended this result to systems reconstructed from the time series of a single variable and test the results for three simple, autonomous nonlinear systems.

The appearance of high growth rate bred vectors proves to be a skillful predictor of regime change within the Lorenz system. Additionally, the number of high growth rate bred vectors in a given regime is correlated with the duration of

Table 3.6: Measures of forecast accuracy in terms of the Hit Rate (HR), Threat Score (TS), and False Alarm Rate (FAR); (a)-(c) are the same as in Fig. 3.18. The final row shows the values when the threshold of $x(t_i)$ rule is used.

	PC (%)	HR (%)	TS (%)	FAR (%)
(a)	3.7	1.0	0.7	89.6
(b)	3.6	2.1	1.5	92.7
(c)	65.7	70.4	59.3	45.6

the subsequent regime, allowing for prediction of the number of orbits in the next regime.

The Chua attractor also has two scroll regimes. The mechanism of regime change in this system differs from that of the Lorenz system. High growth rate bred vectors tend to appear prior to the regime change, however their skillfulness as a precursor for regime change is not as strong. In the phase space reconstructed from a single model variable, using high growth rate bred vectors to forecast the regime change is more accurate.

The Rössler attractor has a structure that is simple compared to either the Lorenz or Chua attractors. When using knowledge of the full phase space, either through the equations of motion or a time series of all three model variables, the regions in which bred vectors of various growth rates appear are clearly defined. In the case of the reconstructed phase space, those regions become intermingled. However, there is very little predictive skill gained from using the appearance of high growth rate bred vectors to forecast regime change in the case of full knowledge of the phase space. There is some skill when making predictions in the reconstructed phase space using nearest neighbor bred vectors.

The ability to predict regime change in a dynamical system using the time series data of just one of its many variables, demonstrated here, has important implications. For most systems in nature and in laboratory, the time series observations of only a limited number of physical variables, often a single variable, are available. In many cases even the actual number of variables is not known. The results presented here demonstrate that the nearest-neighbor breeding enables the prediction of regime change in systems for which regime change follows the appearance of instabilities, thus extending the predictive capability beyond the cases whose time evolution equations are known. Further, when regime change is associated with large changes in the dynamical states, this technique can lead to the prediction of large or extreme events in the cases where nonlinear dynamical predictions are made using time series data, e.g., in the magnetosphere and space weather as will be seen in Ch. 5.

Chapter 4: Modeling and Forecasting the Lorenz System

4.1 Introduction

We have seen that the predictability of the Lorenz system is preserved in the phase space reconstructed from a single variable, and that a phase space model constructed by the method of time delay vectors can serve as a model for making forecasts. Now we wish to demonstrate that data assimilation techniques, in particular the Ensemble Transform Kalman Filter (ETKF), can be applied to forecasts made using data driven models rather than the more typical application to numerical models.

As a test bed for determining whether the NN ETKF method described in [2.5](#) serves to improve forecasts, we again turn to the Lorenz 63 system. The dynamical equations are well known and easy to integrate to forecast the model forward in time. The reconstruction of the Lorenz system is described in [3.2.1](#).

When forecasting the future state of the Lorenz system using the true dynamical equations, the model is perfect despite any error in the observations. In the case of a perfect phase space model, it is possible to directly compare the results of using the NN ETKF in either the full phase space of the model or the phase space reconstructed from a scalar time series of a single model variable to the ETKF method

using the model equations to make forecasts.

The goal of this work is to demonstrate that data assimilation techniques can be used to improve forecasts made of real world observables for which modeling numerically is not possible. In these cases, historical observations from which data driven models would be constructed are contaminated by observational noise. Thus, the phase space model constructed from the noisy time series will also contain errors.

Therefore, we will study the effects of random errors on the quality of the forecasts. In real-world systems, we are not so fortunate as to know the underlying equations and observations are tainted by instrument error. The ability of NN ETKF to adapt to these noise measurements will be important for assessing its applicability to real world problems. To remove some of the noise from the data, we apply Singular Spectrum Analysis as described in [2.3](#). Noise can enter the system in many ways. The most straight forward is instrument noise that is added to observations of the system variables. However there may also be parameters that vary within the system. To simulate model errors, a series of random perturbations will be added to the parameters of the Lorenz system in [Eq. 3.2](#).

4.2 The Perfect Model Case

A “perfect model” is one that introduces no error to the forecasts. In the case of a traditional ETKF scheme with a numerical model, this means using the known dynamical equations, [Eqs. 3.2](#), for the Lorenz attractor to make predictions for the evolution of points on the attractor. For data derived models, this means the

model is constructed from the exact solutions of the dynamical equations without any additional noise or error added.

The benefit of the perfect model is it allows for direct comparison between the forecasts made by solving the dynamical ODEs from the exact analysis initial conditions to data based approaches. As in Ch. 3, three experiments are conducted. Experiment (a) is the standard application of the ETKF using Eq. 3.2 to forecast. Experiment (b) is the application of the NN ETKF to a phase space model that includes all of the phase space variables (x, y, z) using a long time series obtained by integrating Eqs. 3.2. From this long time series, we can take the first variable $x(t)$ and construct an embedded phase space model to test the method developed in this study and outlined in Sec. 2.5. The third experiment, (c), takes this reconstructed phase space model with components (x_1, x_2, x_3) for an $m = 3$ dimensional embedding with $\tau = 7$ to make forecasts using the NN ETKF.

The data-derived models are constructed from the same set of solutions to the Lorenz ODEs containing $N = 200\,000$ points. This means that the first coordinate of the full phase space model is the same as the coordinate in the reconstructed phase space model with the latest time index, i.e. x_3 . To compare the three methods, each will be used to predict the value of a predetermined control trajectory, also obtained by solving Eq. 3.2. The control trajectory is out of sample from the training data used to construct the model for the data driven techniques. Observations will be taken from this truth trajectory and also are perfect in the sense that no observational error has been added. To ensure that the same observations are used in each of the three experiments, the first variable of the observation vector will be

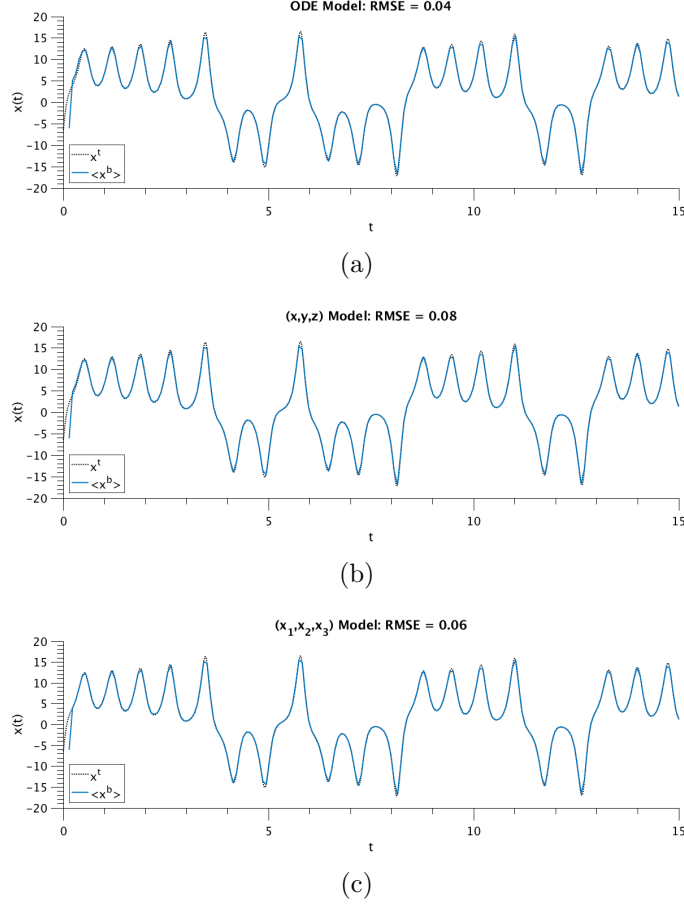


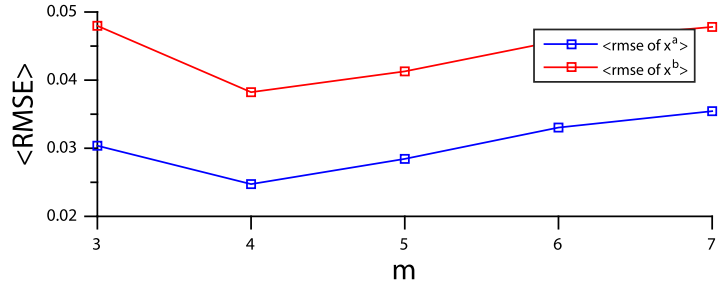
Figure 4.1: Comparison of forecasts made using three versions of the perfect model: 4.1(a) the ODE model, 4.1(b) the full phase space model (x, y, z) , and 4.1(c) the phase space model (x_1, x_2, x_3) reconstructed from $x(t)$.

observed for the ODE and full phase space experiments, and the third variable x_3 will be observed in the embedded phase space experiment. In each case, the scalar variable observed for each forecast window will be the same. The initial ensemble is selected from randomly chosen points on the full phase space attractor.

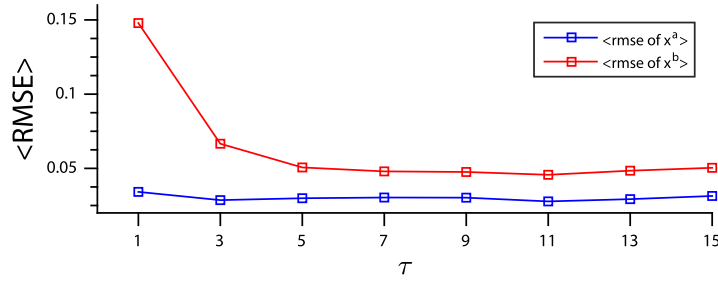
4.2.1 Results for the Lorenz System

Results are shown for forecasts of 8 time steps in Fig. 4.1. To compare the forecasts, the results for forecasting the x variable are shown. In the case of

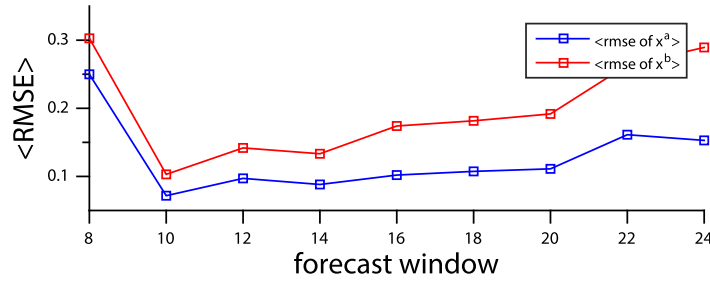
the reconstructed phase space model, this is obtained following the reconstruction procedure outlined in Sec. 2.5. The black curve in each panel represents the true state of the system the forecast attempts to capture and is the same for all. The colored curves represent the analysis and background determined using ETKF the various forecasting techniques, as described in the caption.



(a)



(b)



(c)

Figure 4.2: The average root mean squared error in both the analysis (blue) and the background (red) as a (a) function of embedding dimension, (b) time delay, and (c) the length of the forecast window.

Each of the three methods produces both forecasts (the darker colored curve)

and analysis (the lighter colored curve) values that are in close agreement with the truth. The RMSE errors in the background forecast for each of the three techniques are also comparable. In the case of using ETKF with the ODE equations and NN ETKF in the full phase space, knowledge of the full information of the system is used. The forecasts are excellent, as would be expected when the observations also contain no error. The remarkable thing here is that when only a time series of a single variable is known, the same results can be recovered.

4.2.2 Tuning Embedding Parameters

Another benefit to the perfect model is that the effect of various parameters on the nearest neighbor ETKF forecast can be compared. Such parameters include the time delay τ , the embedding dimension m , and the duration of the forecast window t_w . Figure 4.2 shows the root mean squared error (RMSE) in each the background (red curve) and the analysis (blue curve) as each of the quantities is varied. In the first panel the time delay is varied while the embedding dimension and forecast window are kept constant at $m = 3$ and $t_w = 8$. In the second panel the embedding dimension is varied and $\tau = 7$ and $t_w = 8$. In the final panel, the forecast duration is increased for a reconstructed attractor with $\tau = 7$ and $m = 3$.

The dimension of the Lorenz attractor is well known and the false nearest neighbors test revealed that the number of dimensions sufficient to unfold the attractor in the reconstructed phase space is 3. Including more dimensions than necessary in the reconstruction adds additional information through the inclusion of

more time delayed values of the x -coordinate. This additional information, however, does not improve the forecast. With additional dimensions in the reconstruction the RMS error in the forecast and analysis both increase as can be seen in Fig. 4.2(a).

A range of values of the time delay produce forecasts of comparable skill given an embedding dimension of 3. The average mutual information function also indicated that this same range would produce a good reconstruction, as discussed in section 2.2. Time delays beyond 15 time steps are too long and the components of the time delay vectors begin to become uncorrelated.

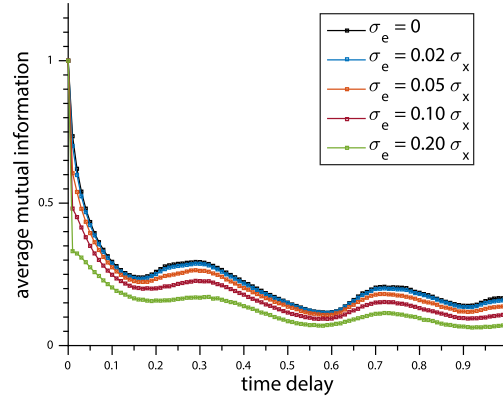
The experiments in this section were conducted using a forecast window of 8 time steps. As the forecast window increases, the skill of the forecast decreases. This is also the case when using the ODE equations to make forecasts.

4.3 Additive Observational Noise

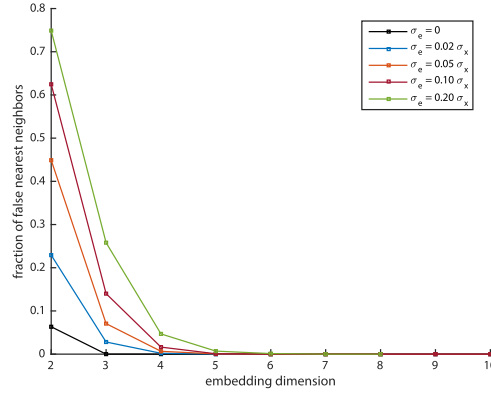
4.3.1 Phase Space Reconstruction of Noisy Data

To simulate the realistic conditions in which observations have noise in the Lorenz model testbed, additive noise is incorporated into both the time series used to construct the model and the observations used during the ETKF cycles. A time series of normally distributed perturbations are added to model trajectory and control trajectory to generate observations. The noise level is then scaled to determine the effect of decreasing signal to noise ratio.

The effect of the noise on the reconstruction is examined. The addition of noise increases the fraction of false nearest neighbors even at higher dimensions. Random



(a)



(b)

Figure 4.3: (a) The mutual information function as a function of time delay and (b) the fraction of false nearest neighbors as a function of embedding dimension for the Lorenz $x(t)$ plus various levels of additive noise variable.

noise is high dimensional and effectively increases the dimension of the attractor.

From Fig. 4.3(b) it can be seen that while the fraction of false nearest neighbors falls to zero after 3 dimensions for the perfect model, the minimum embedding dimension is higher as the noise level increases.

To effectively unfold the noisy attractor, an embedding dimension larger than the minimum dimension indicated by the false nearest neighbors test is used for the time-delay embedding. In this case, 10 dimensional time-delay vectors are con-

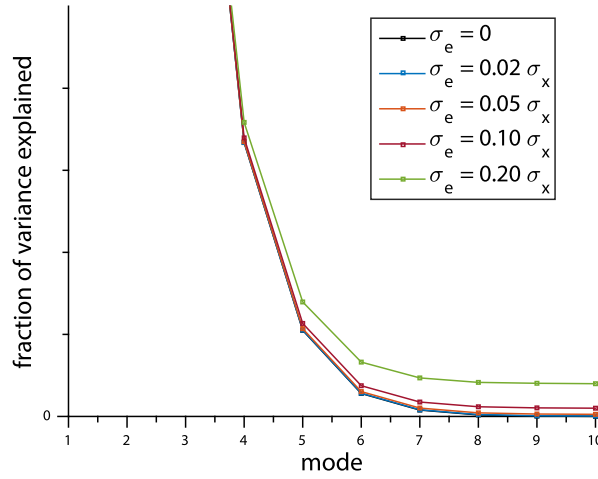


Figure 4.4: An enlarged portion of the eigenvalue spectrum computed for various levels of additive noise for the Lorenz $x(t)$ variable to show the noise floor.

structed from the noisy $x(t)$ time series. Again, a time delay of 7 time steps is used since the behavior of the average mutual information function is not affected.

We are interested in forecasting only the deterministic signal even though it is contaminated by noise. Singular spectrum analysis (SSA), described in Sec. 2.3, can be applied to the time delayed vectors in order to separate the modes of variability corresponding to the signal from those of the noise. Eigenvectors are computed from the covariance matrix of the model trajectory covering the attractor. Regardless of the level of the noise present in the signal, the eigenvectors are roughly the same for each noisy trajectory, indicating that they correspond to the deterministic part of the signal. The first six eigenvectors for a model trajectory with no error and an error level of 50% of the variance in the signal are shown in Fig. 4.5. The curves are essentially identical.

The spectrum of eigenvalues shows the noise floor where the value tends to level off. From the spectrum of eigenvalues, it is evident in Fig. 4.4 that with

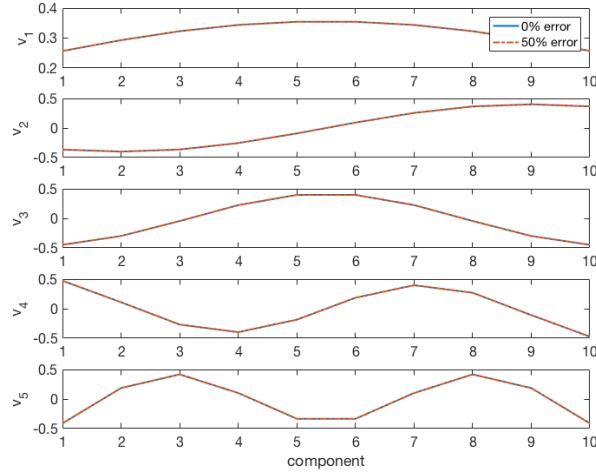


Figure 4.5: The first six eigenvectors of the Lorenz system x variable with additive noise levels of 0% and 50%.

increasing magnitude of noise added to the data, the noise floor feature becomes more prominent. However, the noise floor appears roughly after 6 dimensions. The variability represented by those first 6 eigenmodes corresponds to the variability of the deterministic signal (the x -coordinate) in the noisy signal.

Therefore, the first six eigenvectors are retained to serve as the basis for the model. Forecasts of all six phase space variables are made in the model space. The six principal components are then reconstructed as described in Sec. 2.5 to obtain the forecast of the variable $x(t)$.

Figure 4.6 shows the NN ETKF background and forecast for NN ETKF as well as the true value of the variable $x(t)$ without added noise, and the observations of the the $x(t)$ variable with noise added. The NN ETKF forecasts are made in three different scenarios. The first is NN ETKF in the full phase space in Fig. 4.6(a). The second is NN ETKF in the reconstructed phase space in Fig. 4.6(b). And finally, the NN ETKF in the phase space spanned by the first six eigenvectors of the SSA

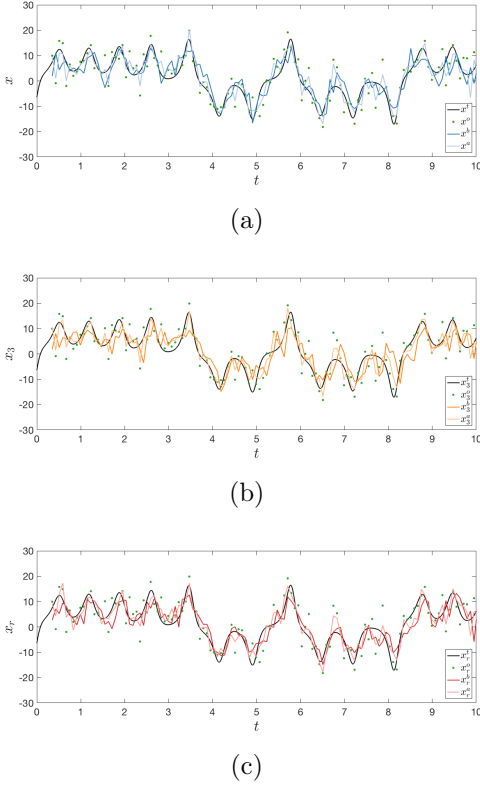


Figure 4.6: Comparison of reconstructed forecasts made of the Lorenz $x(t)$ variable with a noise level 60% of the standard deviation of the x variable. Panel (a) depicts the background and analysis of NN ETKF forecasts made in the full (x, y, z) phase space. Panel (b) depicts the background and analysis of NN ETKF performed in the reconstructed (x_1, x_2, x_3) phase space. Panel (c) depicts the background and analysis of NN ETKF forecasts in the space of the first six principal components after SSA was applied to the reconstructed phase space.

decomposition of the reconstructed model data. In this case, the forecasted values of x come from a reconstruction of the forecasted principal components.

The RMS error in the background and analysis is computed by comparing the mean value of the forecast ensemble x to the true value from the control run. The RMS errors for various noise levels and the three experiments described above are shown in Fig. 4.7. While the best performance for the background comes from NN ETKF in the full phase space, the application of SSA to the reconstructed phase space offers an improvement to forecasts over those made in the reconstructed phase

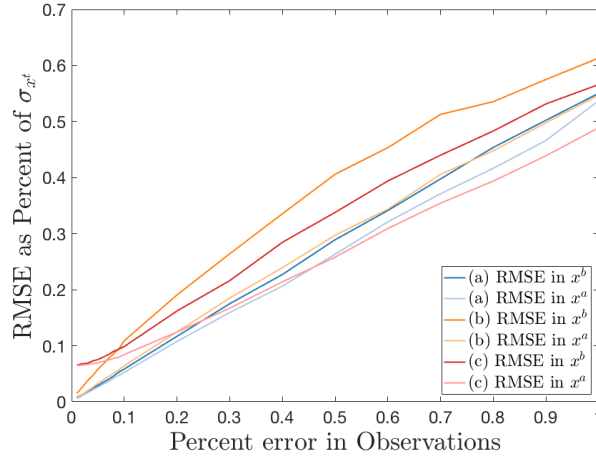


Figure 4.7: The root mean squared error in the mean values of the forecast and analysis ensembles, as a percentage of the standard deviation in the time series $x(t)$ vs. the percentage of error added to the time series to simulate observations.

space alone.

Hamilton et al. [23] also combined phase space reconstruction with ETKF, called the Kalman-Takens filter. This study both confirms their result that the performance of the ETKF in the phase space reconstructed from a single variable is comparable to methods with full knowledge of the phase space or equations, and also extension to this work. The addition of SSA to the construction of the phase space model helps the data driven method to handle the noisy data.

4.4 Discussion and Conclusions

Using only a scalar time series $x(t)$ the dynamics of the Lorenz model can be reconstructed using time-delay embedding techniques. Phase space models constructed from a long time series of observations and NN ETKF are used to forecast future values of the time series $x(t)$. The forecast skill is comparable to the skill

obtained by using numerical solutions of the Lorenz ODEs and ETKF to forecast based on the same initial conditions as seen in Fig. 4.1. Forecasts made using the NN ETKF in the phase space reconstructed from the time series $x(t)$ are also comparably skillful, thus validating this method.

The error in forecasts is not particularly sensitive to the specifics of the parameters used in the embedding procedure. A range of time-delays can be used and similar results obtained. Including additional dimensions beyond the minimum indicated by the false nearest neighbors test does not improve the result. Including many more dimensions than required actually can have a negative impact on the skill of the forecast.

When simulated noise is added to the time series used for phase space reconstruction, forecasts of the time series $x(t)$ can still be made by filtering out the noise using SSA. As the noise level increases, the dimension of the attractor also increases since random noise is high dimensional. Systems with a high level of noise contamination produce less skillful forecasts than perfect model scenarios and low level noise contamination. This indicates that it will be possible to apply these techniques to real data for which contains noise and an unknown signal. The SSA can be used to filter out some of the noise and produce more skillful forecasts than NN ETKF alone.

Chapter 5: Modelling and Forecasting Space Weather Using Time Series Data of Magnetic Field Variations

5.1 Introduction

5.1.1 Geomagnetic Substorm Dynamics

The near-Earth space environment consists of the portion of the atmosphere that extends into the interplanetary medium. At these altitudes the atmosphere becomes ionized as energetic photons strip gas atoms of their electrons. The Earth's magnetic field becomes an important force governing the behavior of the plasma in this region. Both the ionosphere and the magnetosphere interact with the plasma and magnetic field structures carried by the solar wind emanating from the corona of the sun. The term space weather describes the fluctuations in the near-Earth space environment driven by the solar wind. The structure of the sun's magnetic field is complicated and evolves over both short time scales and the eleven year solar cycle. Due to the turbulent nature of the plasma motions within the sun, the solar magnetic field can become quite complex, particularly during the maximal phase of the solar cycle when solar activity is increased.

The earth is shielded from the solar wind by its standing dipole field. The

protective bubble of the magnetosphere segregates the plasma populations within it from mixing with those carried by the solar wind and allows life on earth to exist. The shape of the Earth's dipole field is distorted by the solar wind. On the day-side of the Earth, facing the sun, the dipole field is compressed at the magnetopause, the boundary of the magnetosphere, where the two fields meet. On the night-side the dipole field is elongated by the flow of the solar wind into a structure called the magnetotail. Under certain conditions of the solar wind, the highly energetic plasma carried by the solar wind can enter the magnetosphere and plasmasphere protects, resulting in space weather conditions observed by the Earth. An example of such a condition is the so-called magnetospheric substorms observed at high latitudes.

In the region where the two fields meet at the magnetopause, labeled (2) in Fig. 5.1, the orientation of the Earth's magnetic field is northward. If the magnetic field of the solar wind has a component with a southward orientation it is possible for magnetic reconnection to take place. At the point where the closed line of the Earth and the open line of the solar wind meet, reconnection creates open field lines with one foot point in the Earth's polar regions that opens into space. These open lines are drawn across the polar region to the night-side of the earth. There magnetic tension compresses the lines in the magnetotail region. When the open lines come close enough to one another, reconnection again takes place in region (3). This reforms an open solar wind field line that is carried away by the solar wind and a closed field line for the Earth. The newly formed closed field line is carrying the energetic plasma from the sun. The release of magnetic energy due to reconnection causes the plasma to flow along the closed field line where it precipitates into the Earth's

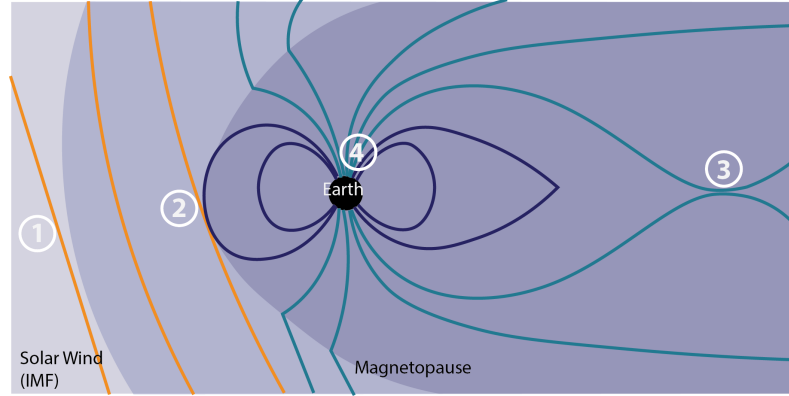


Figure 5.1: Schematic of the interaction between the solar wind magnetic field and the Earth's magnetic field

upper atmosphere over the polar regions. The injection of the highly energetic solar particles into the ionosphere gives rise to the auroras seen in the high latitudes. To close the current loop created by plasma flowing along the field lines, currents are induced in the ionosphere, intensifying the naturally occurring currents found there.

In addition to the visible auroras, the signature of this geomagnetic substorm can be seen in the fluctuations in the Earth's magnetic field caused by the magnetic field induced by the intensified currents in the ionosphere. Ground based magnetometers, located within the auroral oval as shown in Fig. 5.2, record these fluctuations and present the data every minute. To characterize the state of the ionosphere, indices are constructed from the 12 stations located in the auroral oval (Fig. 5.2). The World Data Center at Kyoto University, Japan and other locations worldwide collect data from all the stations (e.g. SuperMAG) and plots anomalies simultaneously by removing a quiet-time reference value from each station. Indices are defined by the upper and lower envelope of the data, called AU and AL respectively shown in Fig. 5.3. The overall measures of state is characterized by

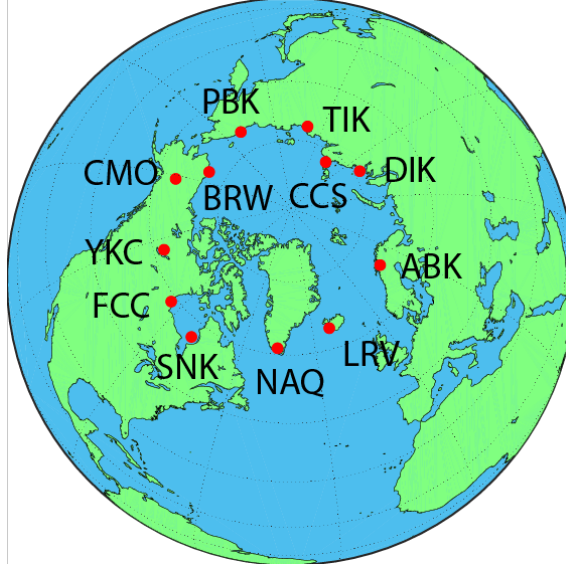


Figure 5.2: The locations of the 12 ground-based magnetometer stations whose measurements contribute to the construction of the AE indices

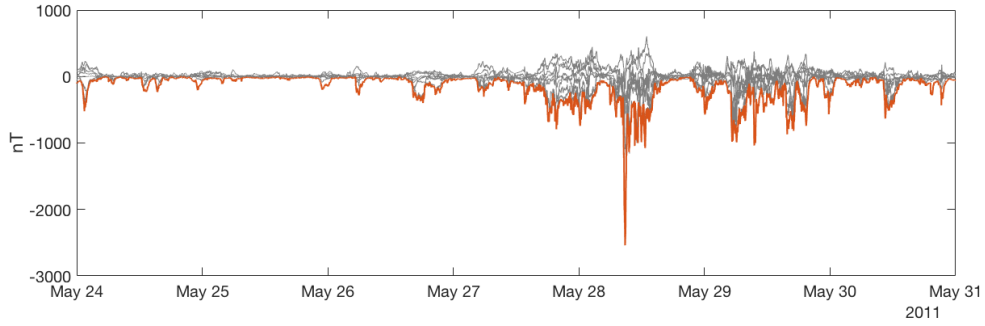


Figure 5.3: The AL index (orange curve) is the lower envelope of anomalies measured at each of the stations (gray curves). Data from eight of the twelve stations that contribute to the AL index are shown during a high speed stream in May of 2011.

$$AE = AU - AL \text{ and } AO = AU + AL.$$

The AL index is almost always negative, as can be seen in Fig 5.3. Since this index corresponds most closely with the intensification of the westward electrojet that takes place during substorm event, it is the index that we will use to forecast these events. Substorms correspond to the large, negative spikes seen in the AL index. Forecasting future values of the index will allow for substorm forecasting.

Forecasting the anomalies measured by individual stations will allow for spatial information about the location of the disturbance and potentially provide information relevant to people within the vicinity of the magnetometer station.

Forecasts will be made using reconstructed phase space models as described in Sec. 2.1. Because the time series of the AL index is derived from observed fluctuations in magnetic field measurements which in turn are driven by the rapidly fluctuating solar wind magnetic fields, the time series data for the reconstruction is contaminated by observational errors. Singular spectrum analysis will be used to isolate the deterministic signal for forecasting. Forecasts will be made using the nearest neighbor ensemble transform Kalman filter described in Sec. 2.4 and applied to the Lorenz system in Ch. 4. Forecasts of the AL index will be made, as well as of the station observations that yield the auroral indices. Additionally, the ensembles used for forecasting these time series will also be studied further to gain insight into the occurrence of extreme events.

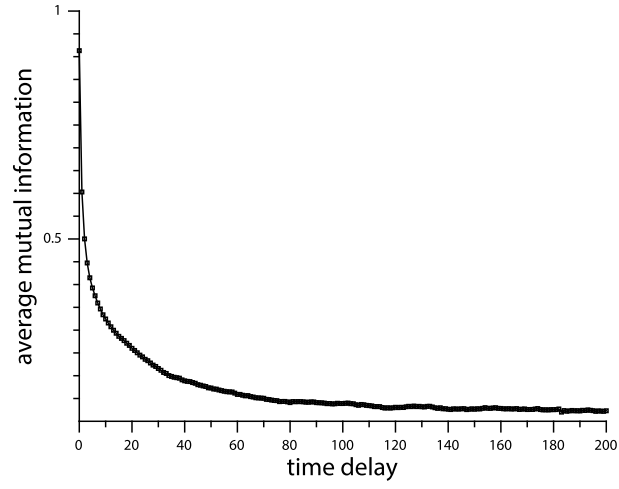
5.1.2 Reconstruction of Magnetospheric Dynamics

The system that yields the auroral electrojet indices consists of many complicated, coupled systems including the solar wind, magnetosphere, and ionosphere. On first inspection it is not obvious that this system is low dimensional, or even deterministic, since the solar wind driver is arguably stochastic in nature. The dimensionality of the AE index, in particular, has been well studied [47–55] Many of the same conclusions can be drawn from the AL index, chosen for analysis due to its

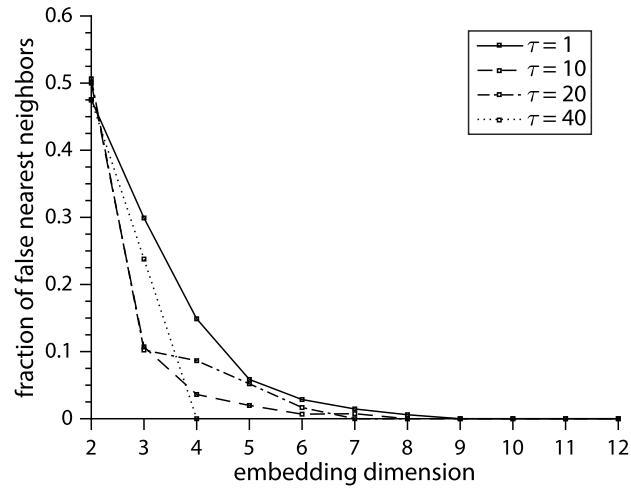
close physical relation with the substorm processes we are interested in forecasting.

The average mutual information of the AL index falls off exponentially with increasing time lag and has no clear first minimum, unlike the same quantity computed for the Lorenz system variable $x(t)$ in Sec. 3.2.1. In such cases, one choice of time delay is a single sampling time step for the reconstruction [24]. Using a time delay of the one 1-minute time step, the fraction of false nearest neighbors can be computed for various embedding dimensions. When the fraction of false nearest neighbors becomes negligible, a sufficient embedding dimension has been reached. All of these suggest that the time series originates from an underlying low dimensional, deterministic system and thus the data assimilation and forecasting technique described above are applicable.

The final step in the analysis of the AL index is to determine the most important modes of variability to isolate the deterministic part of the signal, as was done for the Lorenz attractor with additive noise in Sec. 4.3. The time delay used for the embedding is a single time step, therefore the more important embedding parameter becomes the embedding dimension. Each component of the time delay vector is adjacent in the original scalar time series, thus the time delay vector acts like a sliding window along the time series. This window should be large enough that the attractor is fully unfolded, thus greater than the minimum embedding dimension. Due to the noise contamination seen in the data, the standard deviation of the data over the window should be greater than the noise level so that the signal can be resolved [24]. From $m = 20$ dimensional time-delay vectors, the first six modes of variability are selected. Therefore the dimension of our phase space model is six.



(a)



(b)

Figure 5.4: The mutual information function as a function of time delay and the fraction of false nearest neighbors as a function of embedding dimension for the AL index.

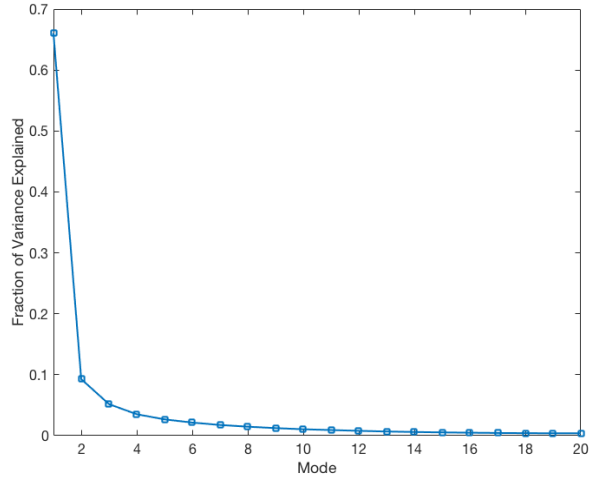


Figure 5.5: The percent variance explained by each mode of the eigenspectrum on the AL index. After 6 modes, nearly 90% of the variance in the signal is explained.

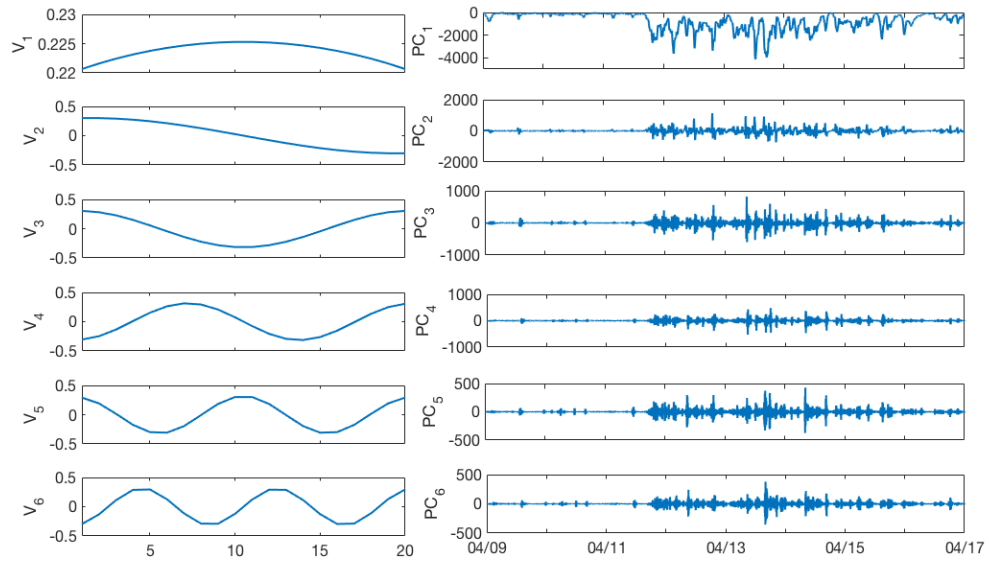


Figure 5.6: The first six eigenvectors and principal components for a substorm event in 2005.

We require a long time series to construct the phase space model for the system. It should include many events of the type we are attempting to forecast so that suitable analogs for the initial conditions of interest can be found. Because the dimension of the underlying attractor is quite high, the duration of the model data must allow for this. Data sets of a year of AL index observations from 2000 are used to construct the phase space models.

Once the model data has been embedded, the principal components are computed as seen in Fig. 5.6. Forecasts will be made in the multidimensional space of the principal components, and then the AL index reconstructed from these forecasted values.

When multivariate data is available, such as station magnetometer measurements, the phase space can be reconstructed in much the same way as for the scalar data. In this case, each variable is physically the same quantity, but obtained at a different location as described in Sec. 2.6. Time-delay vectors for the model are constructed by embedding each channel.

5.2 Forecasting the AL Index During Geomagnetic Substorms

5.2.1 Experimental Setup

We are interested in forecasting magnetospheric substorm events resulting from the interaction between the magnetosphere and the solar wind as described in Sec. 5.1.1. The features of the solar wind that provide favorable conditions for substorm developments are the magnitude and direction of the magnetic field and increases

plasma density and velocity. There must be a persistent southward component of the solar wind magnetic field. The velocity must be high compared to the quiescent solar wind with anomalous plasma density. There are several solar wind conditions that can result in a substorm. - one occurring frequently and one more seldom.

In the following sections, we will look at forecasts of the AL index during a series of substorm events resulting from high speed streams (HSS) of solar wind and coronal mass ejections (CME). The events occur between 2005 and 2014. Most of the events occur during the ascending phase of solar cycle 24, which began in 2008 and reached its peak activity between 2011 and 2014 [56]. To model events during this cycle, data from the year 2000 is used which was near the peak of the previous solar cycle. The model to forecast the AL index was constructed as described in Sec. 5.1.2.

The eigenvectors from SSA performed on the model data form the basis upon which the observations from the events are projected. Forecasts are made in the space spanned by the first six eigenvectors. Ensemble forecasts of 10 members are made. Since the window spanned by the time delay vectors is 20 minutes, forecasts windows longer than 20 minutes should be used to avoid repeating observations during the NN ETKF step. Forecasts of various duration, including 20, 40, and 60 minutes, were made for a variety of substorm events driven by HSS and CMEs.

5.2.2 Forecasting the AL Index during High Speed Stream Events

Generally, the speed of the solar wind is about 350 km/s [57]. The solar wind tends to come from the equatorial region of the sun. Regions of the sun with a single polarity can form, known as coronal holes. While these features can be present at any time, they tend to occur more frequently and persist for longer periods during solar minima.

The solar wind emanating from these holes tends to have higher speeds, often exceeding 600 km/s [57]. During a high speed stream (HSS) event, a high velocity solar wind is emitted, which compresses the lower velocity solar wind ahead of it resulting in co-rotating interaction regions depicted in Fig. 5.7. The increased density of the plasma and magnetic field lines ahead of the high speed stream due to interactions with the lower speed background solar wind produces a solar wind event that has both high velocity, larger plasma density, and increased field strength, thus increasingly providing the proper conditions for a geomagnetic substorm event.

Table 5.1 contains the dates of the HSS events in this study. They occur throughout the period of interest and during both the solar minimum and maximum of solar cycle 24. The geomagnetic substorms that occurred during these events also cover a range of intensities. During the weakest substorm during May of 2011, the AL index reached a minimum value of only -771 nT. The most extreme substorm saw a minimum AL index value of -2539 during May of 2011. Table 5.1 also reports

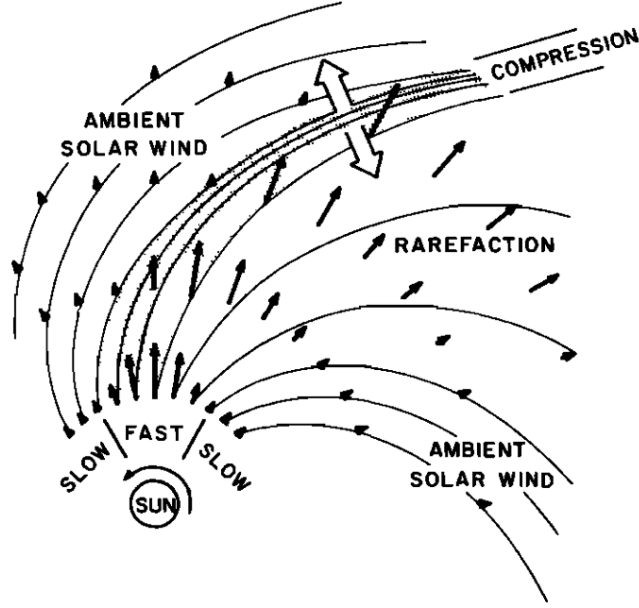


Figure 5.7: Schematic of the co-rotating interaction region (CIR) that forms during an HSS event [1]

the normalized root mean square forecast error given by

$$\text{NRMSE} = \frac{1}{\sigma} \sqrt{\frac{1}{N} \sum_{i=1}^N (x_i^o - x_i^b)^2}, \quad (5.1)$$

where σ is the standard deviation of the observed AL index, x_i^o is the observed value of the AL index, and x_i^b is the forecasted value either using the NN ETKF or persistence.

In general, the forecasts made using NN ETKF show an improvement over those made using persistence evident by reduced normalized RMS errors. This improvement tends to be greater the longer the forecast duration. An example of forecasts made of a single event, that from April 2011, is shown in Fig. 5.8. Skill scores defined by Eq. 2.15 are quoted above each forecast. The skill score is

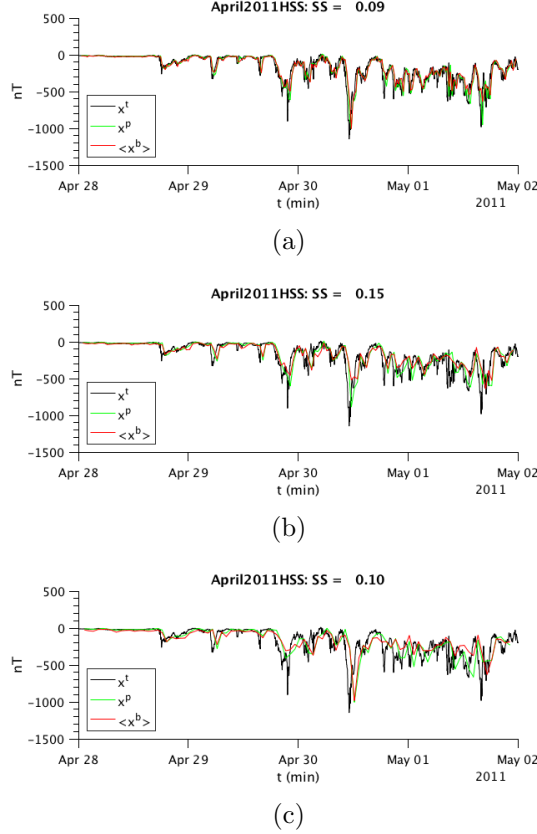


Figure 5.8: Forecasts made of the AL index during an HSS event that occurred in April of 2011. The green and red curves represent forecasts made using persistence and NN ETKF respectively. 5.8(a) depicts a series of 20 minute forecasts, 5.8(b) depicts a series of 40 minute forecasts, and 5.8(c) depicts a series of 60 minute forecasts. The skill score of the NN ETKF forecasts with respect to persistence is quoted in the title.

positive for each case indicating that the NN ETKF forecasts improve the error in the forecast with respect to persistence.

This is a relatively large substorm event with a minimum AL index of -1156. It is difficult for forecasts made using an ensemble of analogs to capture the deepest peaks in AL index because there are relatively few analogs among the model data to use. Thus, as the number of ensemble members increases, more of the members will be forced to use analogs of insufficient intensity resulting in a smoothing out of high intensities in the mean of the forecast ensemble.

Table 5.1: Normalized root mean squared errors for forecasts of various HSS events. Forecasts are made using NN ETKF and persistence, predicting the value of the AL index 20, 40, and 60 minutes beyond the current observation.

	20 min. forecast		40 min. forecast		60 min. forecast	
HSS Event Date	NN ETKF	Persist.	NN ETKF	Persist.	NN ETKF	Persist.
Apr. 9-17, 2005	0.5355	0.5444	0.7036	0.7473	0.7574	0.8655
Apr. 6-12, 2006	0.6446	0.6445	0.7607	0.7753	0.8836	0.9276
Jul. 25-30, 2006	0.6191	0.6528	0.7425	0.798	0.8787	0.9841
Nov. 7-13, 2006	0.4711	0.513	0.6003	0.6899	0.7871	0.9476
Dec. 3-9, 2006	0.4889	0.5158	0.6941	0.7460	0.7078	0.7775
Jan. 28-Feb. 1, 2007	0.6268	0.6614	0.8291	0.8775	0.9805	1.0853
Oct. 23-28, 2007	0.5294	0.549	0.7033	0.7676	0.7924	0.8862
Mar. 6-17, 2008	0.5326	0.5476	0.7453	0.8164	0.8156	0.9132
Apr. 19-28, 2008	0.5089	0.5396	0.7385	0.7585	0.8324	0.9205
Sept.1-10, 2008	0.4761	0.4815	0.6626	0.6942	0.6831	0.7445
Oct. 20-28, 2010	0.6231	0.6385	0.8546	0.9064	0.9135	1.0304
Feb. 26-Mar. 4, 2011	0.5415	0.5455	0.6998	0.7740	0.7652	0.8299
Apr. 28-May 2, 2011	0.5388	0.5660	0.6785	0.7346	0.8379	0.8837
May 24-31, 2011	0.5765	0.6086	0.7260	0.7876	0.7409	0.7887
Apr. 8-15, 2012	0.4848	0.4748	0.6747	0.7059	0.7754	0.8000
May 7-11, 2012	0.5519	0.5513	0.6944	0.7844	0.7242	0.8045
Jun. 28-Jul. 3, 2012	0.5643	0.5762	0.7228	0.7554	0.8584	0.9225
Feb. 26-Mar. 4, 2013	0.4903	0.5186	0.7158	0.7566	0.8487	0.9525

Despite the difficulty in capturing every peak, the forecasts made using NN ETKF tend to have better correlations with the true value of the AL index than forecasts made using persistence. Fig. 5.9 shows the correlation of forecasts with increasing lead times with the true value. Consistently, the correlation coefficient of NN ETKF forecasts is higher than that of persistence.

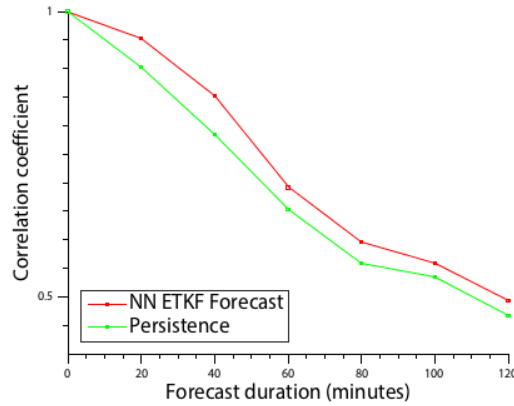


Figure 5.9: The correlation between the true value of the AL Index during the April 2011 HSS event and forecasts made using the NN ETKF and persistence.

5.2.3 Forecasting the AL Index during Coronal Mass Ejections

HSS events tend to occur during solar minima when coronal holes are more prevalent. During solar maxima, the magnetic field of the sun becomes increasingly complex. Active regions of the sun, i.e. where solar magnetic field lines have concentrated, produce magnetic field structures as the layers of the sun are churning causing twisting of the field lines. The exact series of events that cause a build up of solar plasma and twisted magnetic field structures to be ejected from the solar atmosphere are not well understood. However, coronal mass ejections (CME), depicted schematically in Fig. 5.10, tend to form in the coronal streamer belt where closed loop field lines prevent the plasma from escaping into the solar wind [57]. Instead, bubbles build restrained by increasingly unstable magnetic field structures. Eventually these structures are ejected by the release of magnetic potential energy through magnetic reconnection. While the bubble within the CME, seen in Fig. 5.10, contains relatively low density plasma, the propagation of the CME through

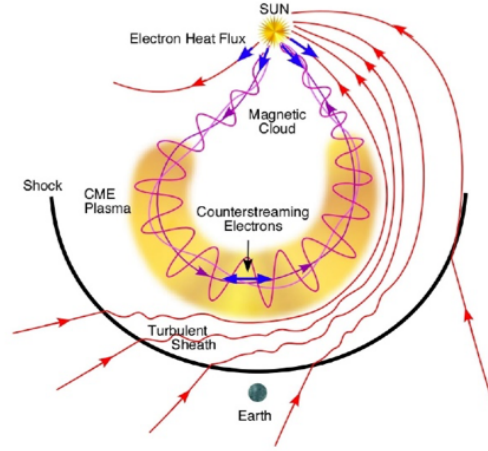


Figure 5.10: Schematic the structure of a CME [2]

the solar wind can form a shock with compressed solar wind plasma ahead of the bubble. Additionally, the flux rope within the CME can have a high magnitude, rapidly varying magnetic field with a strong southward component. CMEs can travel at any speed. Some are carried by the solar wind and some have very high speeds exceeding 2,000 km/s.

Table 5.2 lists the CME events in this study. All occur between 2011 and 2014, during the maximum period of solar cycle 24. The CMEs resulted in substorms covering a range of intensities from the least intense substorm with a minimum AL of -861 in September of 2014, to the most intense substorm occurring in August of 2011 with a minimum AL of -2084. The normalized RMS error in the NN ETKF forecasts and the persistence forecasts is also shown in Table 5.2 for forecasts of 20, 40, and 60 minutes. For shorter duration forecasts of 20 minutes, the NN ETKF method offers little, if any, improvement over persistence. However, for forecasts with a longer lead time, NN ETKF produces a more skillful forecast.

Fig. 5.11 shows forecasts made of a substorm resulting from a CME in Septem-

Table 5.2: Normalized root mean squared errors for forecasts of various CME events. Forecasts are made using NN ETKF and persistence, predicting the value of the AL index 20, 40, and 60 minutes beyond the current observation.

	20 min. forecast		40 min. forecast		60 min. forecast	
CME Event Date	NN ETKF	Persist.	NN ETKF	Persist.	NN ETKF	Persist.
Aug. 2-8, 2011	0.5156	0.4867	0.6964	0.7033	0.8729	1.0525
Sept. 23-30, 2011	0.5673	0.5931	0.7254	0.7684	0.8248	0.9173
Oct. 22-27, 2011	0.4628	0.4560	0.5617	0.5821	0.7741	0.7744
Apr. 19-27, 2012	0.4671	0.4595	0.6013	0.6139	0.6538	0.7504
July 12-18, 2012	0.499	0.5154	0.5690	0.5622	0.6577	0.6894
Sept. 28-Oct. 3, 2012	0.4151	0.4635	0.5499	0.5091	0.5967	0.5953
Nov. 10-16, 2012	0.3827	0.3748	0.5171	0.5311	0.6206	0.5826
Feb. 25-Mar. 2, 2014	0.4204	0.4396	0.5483	0.5696	0.6628	0.68
Apr. 9-15, 2014	0.4531	0.4626	0.6199	0.6611	0.7483	0.7877
Sept. 27-Oct. 2, 2014	0.6108	0.6526	0.8092	0.9206	0.8999	1.0404

ber of 2011. With an even deeper drop in AL during the most intense part of the substorm than that for the HSS event in Fig. 5.8, it is even more evident how the NN ETKF forecast can struggle to match the deep drops in AL. persistence forecasts can produce the feature, however it necessarily lags that of the observed substorm.

These results also suggest that the NN ETKF offers an improvement over other data driven techniques for forecasting substorms via auroral electrojet indices. Chen et al. [58] use a database of AL index from 2001 and solar wind $-B_z v_x$ to reconstruct a phase space model and the local linear filter to forecast substorms from the peak of solar cycle 23 [58]. The largest substorm due to a CME in this study occurred in May of 2011. It was comparable in magnitude to the November 2003 substorm in Chen’s study. Conducting 5 minute forecasts, a normalized RMS

error of 0.792 was achieved. Using the NN ETKF, a 60 minute forecast has a normalized RMS error of 0.741 over a similar period.

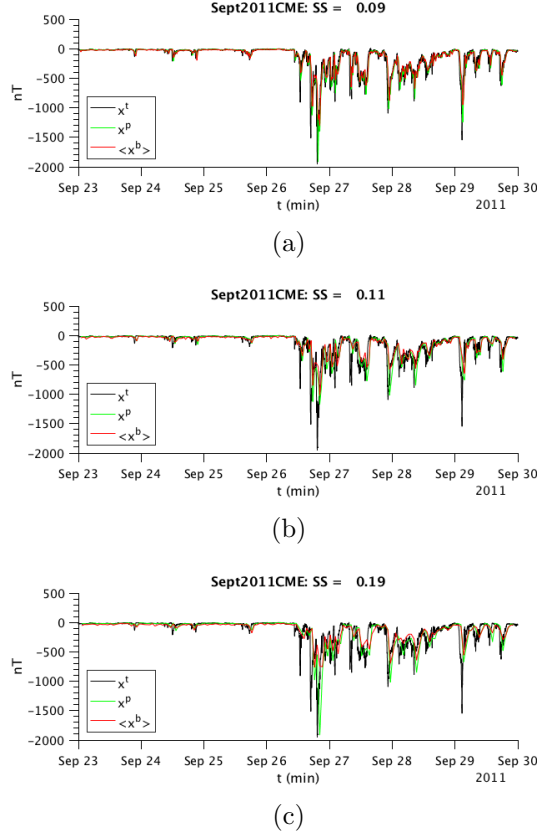
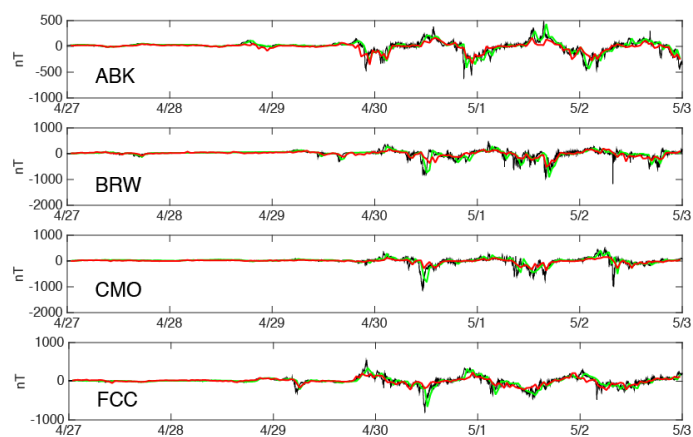


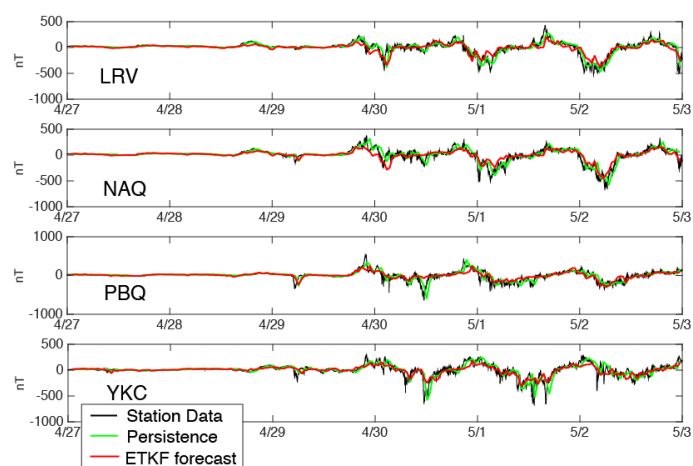
Figure 5.11: Forecasts made of the AL index during a CME event that occurred in September of 2011. The green and red curves represent forecasts made using persistence and NN ETKF respectively. 5.11(a) depicts a series of 20 minute forecasts, 5.11(b) depicts a series of 40 minute forecasts, and 5.11(c) depicts a series of 60 minute forecasts. The skill score of the NN ETKF forecasts with respect to persistence is quoted in the title.

5.3 Forecasting Ground Based Magnetometer Measurements

The magnetometers monitoring fluctuations in the Earth's magnetic field are maintained by various national and multinational organizations. Of the 12 stations located in the auroral oval, the four stations located in Russia do not make their data



(a)



(b)

Figure 5.12: Forty minute forecasts of the data from each of the 8 stations that had data available to the public and the skill scores of the forecasts with respect to persistence. Skill scores range from -0.05 to 0.27.

available to the public and recent data are not available. Therefore, only eight time series were available for a sufficiently long time to make forecasts. The horizontal component, that is the component in the plane tangent to the station's location, were collected from each of these eight stations.

Because the magnetic latitude and longitude locations of the stations vary, each has a different baseline reading on a quiet day with little to no magnetospheric activity. To account for this, averages over the year were removed from each station's reading. The same embedding parameters as for the AL index were used to construct time delay vectors. As described in Sec. 2.6, the time delay vectors are of dimension $m \times L$ where in this case $m = 20$ and $L = 8$ for the eight stations providing data. The model was constructed from data recorded in the year 2000.

The multi-channel SSA extension described in Sec. 2.6 was applied to the trajectory of the model data. The first six eigenvectors were retained to form the basis of the phase space model. The projection of the trajectory matrix onto this basis yields a 48 dimensional phase space model, with six degrees of freedom corresponding to this. Twenty minute forecasts of the April 2011 HSS event were made. Forecasts for each station were made simultaneously. The measurements for each station were then reconstructed. The forecasts compared to persistence appear in Fig. 5.12.

5.4 Forecasting Extreme Events using Ensemble Spread

Identifying and predicting extreme events, particularly in space weather, is an area of active research. To identify regions of instability in the phase space of a system that give rise to extreme events, we develop a data driven approach that takes advantage of the forecast ensembles we have constructed. Just as the bred vectors described in Ch. 3 converge to the direction of largest growth locally, the spread of the ensembles will increase when instabilities cause large growth in various directions. The ensemble spread as given by

$$S^b = \sqrt{\text{Tr}(P^b)} \quad (5.2)$$

where P^b is the background error covariance given by Eq. 2.10. The ensemble spread characterizes the uncertainty in the forecast which will grow as an extreme event occurs and predictability is lost [59].

The concept here is readily illustrated using the Lorenz system. In Fig. 5.13, ensembles originating from various locations on the Lorenz attractor are shown. Ensembles that begin in a relatively stable region tend to experience less ensemble spread with repeated integrations, as seen in the left panel. However, ensembles near an unstable part of the attractor, about to experience the extreme event of regime transition as in the right panel, experience a large growth in their spread as they evolve.

During his time as an undergraduate student in the Department of Atmo-

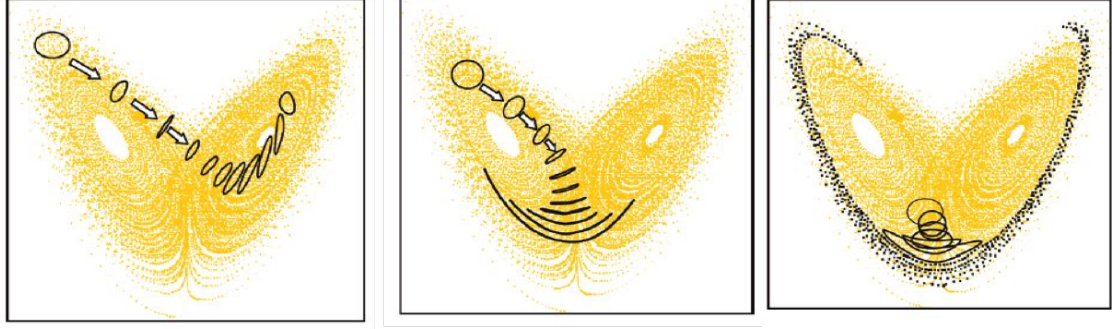


Figure 5.13: The left panel depicts the scenario in which an ensemble is not subject to an extreme event, while the right panel shows the effect of an extreme event on the spread of the ensemble members. Figure from [3].

spheric and Oceanic Science, Keenan Eure and I worked to develop a data driven approach analogous to that developed for bred vectors in Ch. 3 that we test on the Lorenz system. Ensembles about points along a control trajectory will be taken from a densely populated phase space model constructed from a long time series of data. For every 10 time step forecast along the control trajectory, a new ensemble of points within a small neighborhood of the control trajectory is reinitialized. That neighborhood of points is evolved concurrent with the control and the final spread is compared to the initial spread. As the control trajectory approaches an extreme event, in this case the transition from one lobe of the attractor to the other, the spread of the ensemble increases. This is illustrated in Fig. 5.14 by the appearance of large spread ensembles, indicated by red stars, preceding a transition from the variable x from positive values to negative and vice versa.

We extend this to use the ensembles produced using ETKF, which also sample the local phase space estimate the probability distribution of the state space of the model, to forecast values of the AL index. The spread tends to be small in regions where the system is stable, such as during the quiet time prior to the

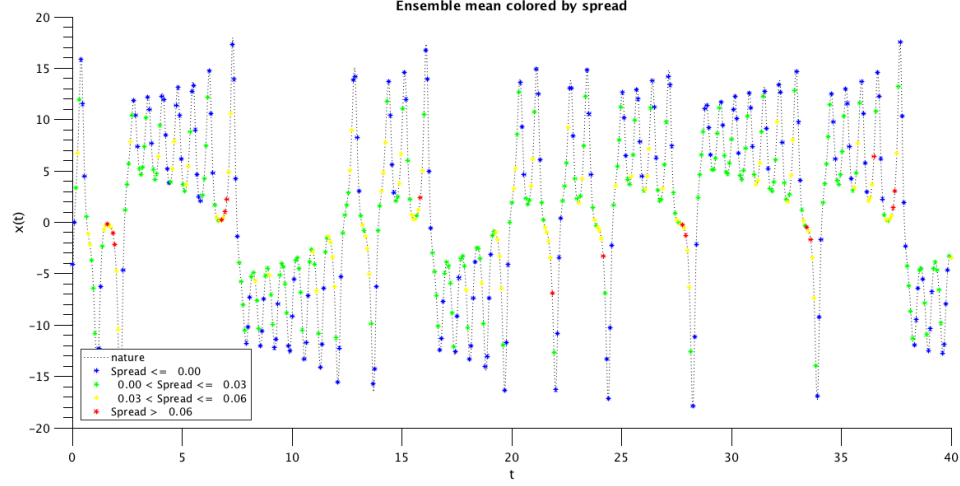


Figure 5.14: The spread of the nearest neighbors to a control trajectory over a series of forecasts is indicated by colored stars along the $x(t)$ variable. This figure was provided by Keenan Eure.

HSS events. As the magnetosphere/ionosphere system becomes more disturbed, the spread increases.

In Fig. 5.15 the time series of the AL index for a HHS event that occurred in April 2011 is plotted along with the spread of the ensemble used to forecast it. The size and color of the dots at the location of each forecast correspond to the size of the ensemble spread. There is a high correlation between the magnitude of the AL index and that of the of the ETKF ensemble spread, indicating that the spread is a good predictor for the magnitude of an event. Events with very large magnitude AL could be considered extreme.

Using the spread and time rate of change of the spread as independent variables, along with the value of the AL index forecasted by the mean of the ensemble generated in the NN ETKF, a prediction of the AL index can be made using linear

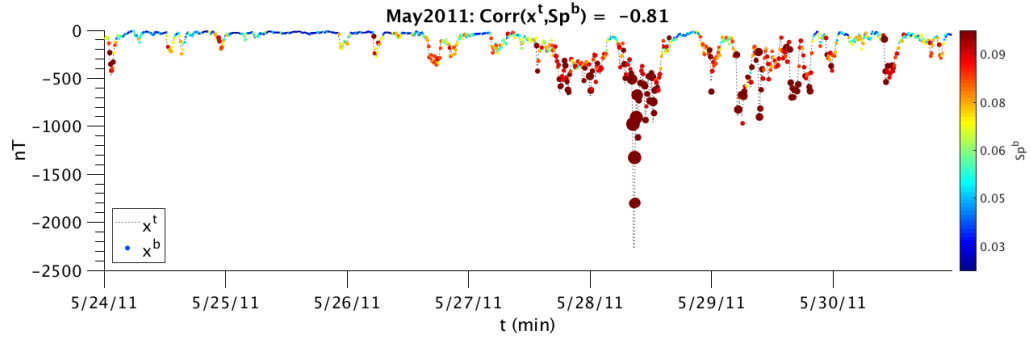


Figure 5.15: The dashed line depicts the observed value of the AL index while the colored circles indicate the mean value of the ETKF forecast ensemble. The color and size of the circles correspond to the magnitude of the ensemble spread. The ensemble spread is well correlated with the value of the AL index.

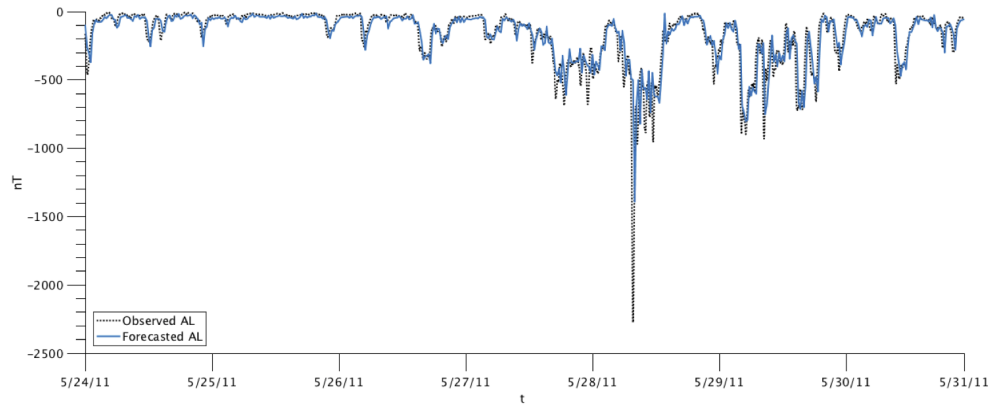


Figure 5.16: The black curve represents the observed value of the AL index during the April 2011 substorm while the blue curve is a forecast made by the linear regression in Eq. 5.3.

regression.

$$AL(t + 1) = b_0 + b_1\bar{x}^b(t) + b_2AL(t - 1) + b_3S^b(t) + b_4S^b(t - 1) \quad (5.3)$$

The estimated value of the AL index from Eq. 5.3 agrees very well with the observed value of the AL index. The training data for the regression coefficients comes from the same model year as used for the phase space models in Sec. 5.2.

5.5 Conclusions

The AL index shows the signature of substorm development and is constructed from anomalies measured by ground based magnetometers in the auroral oval. Non-linear time series analysis of the index indicates that there is an underlying low-dimensional deterministic signal. Reconstruction of the phase space and filtering of high dimensional noise by SSA results in a data-derived model that can be combined with the NN ETKF technique to make forecasts.

Forecasts were made for 18 HSS events occurring during the descending phase of solar cycle 23 and ascending phase of solar cycle 24. Additionally, forecasts for 10 CME events occurring during the peak of solar cycle 24 were also made. While short forecasts do not present a considerable increase in skill with respect to persistence, as the lead time of the forecast increases, so does the skill of the NN ETKF forecasts with respect to persistence. Since persistence forecasts lag the observed values of the AL index, NN ETKF forecasts also demonstrate better correlation with the true value of the AL index.

The index removes all spatial information about the location of disturbances seen. However using the multi-channel extension to the NN ETKF technique allows all of the available station data to be forecast simultaneously. While these forecasts are not always as skillful, proving this location information is valuable. There are local effects of geomagnetic activity that can only be anticipated if regional forecasts are available.

Finally, the spread of the ensembles produced in the NN ETKF forecast offers additional information about the severity of the substorm. As the intensity of the storm increases, fewer analogs provide an appropriate representation of the event. Thus the spread increases rapidly as ensemble members struggle to capture the lowest values of the AL index.

Chapter 6: Summary and Conclusions

The quest to predict the weather has prompted a deep study of the predictability of nonlinear systems. Techniques to characterize instabilities and mitigate the effect of nonlinearities on forecasts have proven invaluable. Despite its complexity and scale, great effort has yielded numerical models for forecasting the state of the atmosphere. Unfortunately, other nonlinear systems present an even greater challenge. Space weather, for example, presents many challenges to following this route including the vast scale of the domain of interest, the number of coupled systems involved in the phenomena, and the lack of observations of relevant processes. It is desirable, and the aim of this thesis, to demonstrate that some of the techniques common in the field of neutral atmosphere weather forecasting can be applied to systems for which no nonlinear model exists or is easily implemented.

As a first step, the technique of bred vectors was applied to data derived models in Chapter 3. Using three simple, autonomous nonlinear models, the Lorenz system, the Chua circuit, and the Rössler system, reconstructed phase space models were constructed using appropriate time delays and embedding dimensions for the first phase space variable. Each of the attractors for the different systems has a unique topology but a common feature among all is the presence of two distinct regimes.

The growth rate of bred vectors was tested as a predictor for regime change in each of the three systems.

The unique features of each system lead to varying degrees of success using the appearance of high growth rate bred vectors as a precursor for regime change. High growth rate bred vectors proved to be the most skillful in the Lorenz system where the hyperbolic fixed point governs the regime change causing trajectories to diverge in its vicinity. Nearest neighbor bred vectors in the reconstructed phase space had performance comparable to methods for which full knowledge of the phase space is utilized.

The Chua system also presented some promise for predicting regime change using high growth rate bred vectors. In fact, in this case the NN bred vectors in the reconstructed phase space showed more skill in predicting the regime change than their counterparts.

The Rössler system is simple by comparison to the other two. One regime consists of a single orbit in each instance. In this case, the high growth rate bred vectors occur immediately following the regime change for methods with full knowledge of the phase space like traditional bred vectors using the ODE and NN bred vectors in the full phase space. In the reconstructed phase space, bred vectors show modest skill in predicting the regime transition as they are not confined to occur following the regime change.

In Chapter 4, the application of data assimilation techniques to data derived models was tested using the Lorenz system. The technique was verified by applying NN ETKF to the best case scenario - a perfect model. When the model is con-

structed using a control trajectory with no error, the performance of the NN ETKF in the phase space reconstructed from the x variable is identical to that of ETKF using the numerical model.

However, if the goal is to apply these techniques to real-world systems for which modeling numerically is difficult or impossible, this does not provide the best test case. To more accurately assess the applicability of this technique experiments were conducted in which the time series data was contaminated by various levels of observational noise. While the NN ETKF performs well in this scenario, as demonstrated by Hamilton et al, [23], further improvements can be made to the data driven model to improve forecasting the true value of the variable.

To filter out the additive noise, eigenvectors corresponding to the underlying deterministic signal are computed from the covariance matrix of the model trajectory covering the attractor. This process, called SSA, produces a reduced dimensional phase space (from that required to unfold the noisy attractor). Regardless of the level of the noise present in the signal, the eigenvectors are roughly the same for each noisy trajectory, indicating that they correspond to the deterministic part of the signal. Forecasts made in the space of principal components of the the time delay vectors projected onto the reduced basis, perform better than forecasts made in the noise reconstructed phase space alone.

Having verified the technique of NN ETKF, we next applied it to a real-world example in Chapter 5. The time series of the AL index is closely related to the development of substorms during active space weather in the Earth’s magnetosphere and ionosphere. These events are the result of solar wind conditions that depart from

the typical ambient solar wind. The solar wind tends to have increased magnetic field strength in the southward direction, increased velocity, and increased density during high speed streams (HSS) and coronal mass ejections (CMEs).

Various substorms resulting from each of these phenomena were forecasted using NN ETKF in the phase space reconstructed from the time series of the AL index and then reduced in dimension using SSA. As was shown for the Lorenz system, this technique produces skillful forecasts. The forecasts are potentially more skillful with longer lead times than for other methods that may be used to forecast auroral electrojet index as a means of predicting substorms.

This technique was then extended to the multivariate time series of the magnetometer stations from which the AL index was constructed. While this forecast is more challenging due to the difficulty in constructing a good model, the forecasts of individual station measurements still show some skill. It is difficult to construct a good model in this scenario because there are often large gaps in the records from individual stations, while the indices constructed from the stations together lack such gaps. Also, the anomalous measurements at any one station are difficult to capture because those experiencing quiet conditions tend to draw the forecast to a quiet state. As with all data driven modeling that uses analogs to represent the observed state, a lack of suitable analogs for each ensemble member can negatively impact the forecast.

This limitation was addressed by including the spread of the forecast ensemble as a potential independent variable for forecasting. As the system enters unstable state, such as those occurring during a substorm, and particularly for events of

very large magnitude, the spread of forecast ensembles generated using NN ETKF tends to increase dramatically. The high correlation between the magnitude of the background ensemble spread and the value of the AL index suggests that it would be a good predictor for future values of the AL index. Combining with previous and current values of the AL index with previous and current values of the ensemble spread as independent variables, skillful predictions of the future value of the AL index can be made using a simple regression model.

While data driven techniques will always have limits, particularly with respect to the quality and quantity of historical data from which to construct models or training data, this work demonstrates how tools commonly applied to numerical models can be extended to these cases. With a great many systems exhibiting the low dimensional underlying dynamics necessary for this kind of phase space reconstruction, the potential to apply data assimilation to improve forecasts made using data derived models is evident.

Bibliography

- [1] V. Pizzo. A three-dimensional model of corotating streams in the solar wind, 1. theoretical foundations. *J. Geophys. Res.*, 83(A12):5563–55572, 1978.
- [2] Thomas H. Zurbuchen and Ian G. Richardson. In-situ solar wind and magnetic field signatures of interplanetary coronal mass ejections. *Space Science Reviews*, 123(1):31–43, Mar 2006.
- [3] T. N. Palmer and R. Hagedorn, editors. *Predictability of Weather and Climate*. Cambridge University Press, Cambridge, 2006.
- [4] M. Casdagli. Nonlinear prediction of chaotic time series. *Physica D*, 34:335–356, 1989.
- [5] L. A. Smith. Identification and prediction of low dimensional dynamics. *Physica D*, 58:50–78, 1992.
- [6] H. D. I. Abarbanel, R. Brown, J. J. Sidorowich, and L. Sh. Tsimring. The analysis of observed chaotic data in physical systems. *Rev. of Mod. Phys.*, 65(4):1331–1392, 1993.
- [7] J. P. Crutchfield and B. S. McNamara. Equations of motion from a data series. *Complex Systems*, 1:417–452, 1987.
- [8] J. D. Farmer and J. J. Sidorowich. Predicting chaotic time series. *Phys. Rev. Lett.*, 59(8):845–848, August 1987.
- [9] M. Ghil, M. R. Allen, M. D. Dettinger, K. Ide, D. Kondrashov, M. E. Mann, A. W. Robertson, A. Saunders, Y. Tian, F. Varadi, and P. Yiou. Advanced spectral methods for climatic time series. *Rev. of Geophys.*, 40(1):1–1–1–41, 2002.
- [10] T. Berry, D. Giannakis, and J. Harlim. Nonparametric forecasting of low-dimensional dynamical systems. *Phys. Rev. E*, 91:032915–1–032915–7, 2015.

- [11] M. Ghil and P. Malanotte-Rizzoli. Data assimilation in meteorology and oceanography. *Advances in Geophysics*, 33:141–266, 1991.
- [12] P. Bauer, A. Thorpe, and G. Brunet. The quiet revolution of numerical weather prediction. *Nature*, 525(7567):47–55, 09 2015.
- [13] E. Kalnay. *Atmospheric modeling, data assimilation, and predictability*. Cambridge University Press, 2003.
- [14] Zoltan Toth and Eugenia Kalnay. Ensemble forecasting at nmc: The generation of perturbations. *Bull. Amer. Meteor. Soc.*, 74(12):2317–2330, December 1993.
- [15] Zoltan Toth and Eugenia Kalnay. Ensemble forecasting at ncep and the breeding method. *Mon. Weather Rev.*, 125:3297–3319, December 1997.
- [16] E. Evans, N. Bhatti, J. Kinney, L. Pann, M. Pena, S. C. Yang, E. Kalnay, and J. Hansen. Rise undergraduates find that regime changes in lorenz’s model are predictable. *Bull. Am. Met. Soc.*, 85(4):521–524, April 2004.
- [17] R. E. Kalman. A new approach to linear filtering and prediction problems. *Trans. AMSE J. Basic Eng.*, 82(Series D):35–45, 1960.
- [18] A. Lorenc. Analysis of methods for numerical weather prediction. 1986, 112:1177–1194, Q. J. R. Meteorol. Soc.
- [19] O. Talagrand. A study of the dynamics of four-dimensional data assimilation. *Tellus*, 33:43–60, 1981.
- [20] G. Evensen. The ensemble kalman filter: theoretical formulation and practical implementation. *Ocean Dynamics*, 53:343–367, 2003.
- [21] Jeffrey S. Whitaker and Thomas M. Hamill. Ensemble data assimilation without perturbed observations. *Mon. Weather Rev.*, 130:1913–1924, July 2002.
- [22] B. R. Hunt, E. J. Kostelich, and I. Szunyogh. Efficient data assimilation for spatiotemporal chaos: A local ensemble transform kalman filter. *Physica D*, 230:112–126, 2007.
- [23] F. Hamilton, T. Berry, and T. Sauer. Ensemble kalman filtering without a model. *Phys. Rev. X*, 6(011021), 2016.
- [24] H. Kantz and T. Schreiber. *Nonlinear Time Series Analysis*. Cambridge University Press, Cambridge, 2003.
- [25] N. H. Packard, J. P. Crutchfield, J. D. Farmer, and R. S. Shaw. Geometry from a time series. *Phys. Rev. Lett.*, 45(9):712–716, 1980.
- [26] F. Takens. *Dynamical Systems and Turbulence*, chapter Detecting strange attractors in turbulence. Springer, Berlin, 1981.

- [27] A. M. Fraser and H. L. Swinney. Independent coordinates for strange attractors from mutual information. *Phys. Rev. A*, 33(2):1134–1140, 1986.
- [28] J. Timmer, H. Rust, W. Horbelt, and H. U. Voss. Parametric, nonparametric and parametric modelling of a chaotic circuit time series. *Phys. Rev. A*, 274(3-4):123–134, September 2000.
- [29] H. U. Voss, P. Kolodner, M. Abel, and J. Kurths. Amplitude equations from spatiotemporal binary-fluid convection data. *Phys. Rev. Lett.*, 83(17):3422–3425, October 1999.
- [30] J.-P. Eckmann and D. Ruelle. Ergodic theory of chaos and strange attractors. *Rev. of Mod. Phys.*, 57(3):617–656, 1985.
- [31] T. Sauer, J. A. Yorke, and M. Casdagli. Embedology. *J. Stat. Phys.*, 65(3/4):579–616, 1991.
- [32] H. D. I. Abarbanel and M. B. Kennel. Local false nearest neighbors and dynamical dimensions from observed chaotic data. *Phys. Rev. E*, 47(5):3057–3068, May 1993.
- [33] M. B. Kennel and H. D. I. Abarbanel. False neighbors and false strands: A reliable minimum embedding dimension algorithm. *Phys. Rev. E*, 66(026209), 2002.
- [34] D. S. Broomhead and G. P. King. Extracting qualitative dynamics from experimental data. *Physica D*, 20:217–236, 1986.
- [35] D.S. Wilks. *Statistical Methods in the Atmospheric Sciences: An Introduction. International Geophysical Series*, volume 59. Academic Press, 1995.
- [36] E. N. Lorenz. Deterministic nonperiodic flow. *J. Atmos. Sci.*, 20:130–141, March 1963.
- [37] E. Lynch, D. Kaufman, A. S. Sharma, E. Kalnay, and K. Ide. Brief communication: Breeding vectors in the phase space reconstructed from the time series data. *Nonlin. Processes Geophys.*, 23:137–141, 2016.
- [38] R. S. Yadav, S. Dwivedi, and A. K. Mittal. Prediction rules for regime changes and length in a new regime for the lorenz model. *J. Atmos. Sci.*, 62:2316, 2005.
- [39] M. J. Hoffman, Eugenia Kalnay, J. A. Carton, and Shu-Chih Yang. Use of breeding to detect and explain instabilities in the global ocean. *Geophys. Res. Lett.*, 36(L12608), 2009.
- [40] A. Norwood, E. Kalnay, K. Ide, S.-C. Yang, and C. Wolfe. Lyapunov, singular and bred vectors in a multi-scale system: an empirical exploration of vectors related to instabilities. *J. Phys. A: Math. Theor.*, 46(254021), 2013.

- [41] D. Vassiliadis, A. Surjalal Sharma, and K. Papadopoulos. Lyapunov exponent of magnetospheric activity from al time series. *Geophys. Res. Lett.*, 18(8):1643–1646, August 1991.
- [42] Leon Chua. The genesis of chua’s circuit. *AEU. Archiv fur Elektronik und Ubertragungstechnik*, 46, 07 1992.
- [43] L. Chua, M. Komuro, and T. Matsumoto. The double scroll family. *IEEE Transactions on Circuits and Systems*, 33(11):1072–1118, November 1986.
- [44] O. E. Rossler. An equation for continuous chaos. *Phys. Lett.*, 57A(5):397–398, July 1976.
- [45] O. E. Rossler. An equation for hyperchaos. *Phys. Lett.*, 71A(2,3):155–157, April 1979.
- [46] A.M. Jade, B. Srikanth, V.K. Jayaraman, B.D. Kulkarni, J.P. Jog, and L. Priya. Feature extraction and denoising using kernel pca. *Chemical Engineering Science*, 58(19):4441 – 4448, 2003.
- [47] D. A. Roberts. Is there a strange attractor in the magnetosphere? *J. Geophys. Res.*, 96(A9):16,031–16,046, September 1991.
- [48] L.-H. Shan, C. K. Goertz, and R. A. Smith. On the embedding-dimension analysis of ae and al time series. *Geophys. Res. Lett.*, 18(8):1647–1650, August 1991.
- [49] A. S. Sharma, D. Vassiliadis, and K. Papadopoulos. Reconstruction of low-dimensional magnetospheric dynamics by singular spectrum analysis. *Geophys. Res. Lett.*, 20(5):335–338, March 1993.
- [50] A. S. Sharma. Assessing the magnetosphere’s nonlinear behavior: Its dimension is low, its predictability, high (us national report to iugg, 1991- 1994). *Rev. Geophys. Supple.*, 33:645–650, 1995.
- [51] J. A. Valdivia, A. S. Sharma, and K. Papadopoulos. Prediction of magnetostorms by nonlinear models. *Geophys. Res. Lett.*, 23(21):2899–2902, 1996.
- [52] D. Vassiliadis, A. J. Klimas, D. N. Baker, and D. A. Roberts. A description of the solar wind-magnetospher coupling based on nonlinear filters. *J. Geophys. Res.*, 100, 1995.
- [53] D. Vassiliadis, A. S. Sharma, T. E. Eastman, and K. Papadopoulos. Low dimensional chaos in magnetospheric activity from time-series ae data. *Geophys. Res. Lett.*, 17:1841, 1990.
- [54] J. Chen. *Spatio-Temporal Dynamics of the Magnetosphere during Geospace Storms*. PhD thesis, University of Maryland, College Park, College Park, MD, USA 20742, 2007.

- [55] S. P. Devi, S. B. Singh, and A. S. Sharma. Deterministic dynamics of the magnetosphere. *Nonlinear Proc. Geophys.*, 20:11–18, 2013.
- [56] Kakad, Bharati, Kakad, Amar, Ramesh, Durbha Sai, and Lakhina, Gurbax S. Diminishing activity of recent solar cycles (22-24) and their impact on geospace. *J. Space Weather Space Clim.*, 9:A1, 2019.
- [57] Tamas I. Gombosi. *Physics of the Space Environment*. Atmospheric and Space Science Series. Cambridge University Press, Cambridge, U.K., 1998.
- [58] J. Chen and A. S. Sharma. Modeling and prediction of the magnetospheric dynamics during intense geospace storms. *J. Geophys. Res.*, 111(A04209), 2006.
- [59] T N Palmer. Predicting uncertainty in forecasts of weather and climate. *Reports on Progress in Physics*, 63(2):71, 2000.